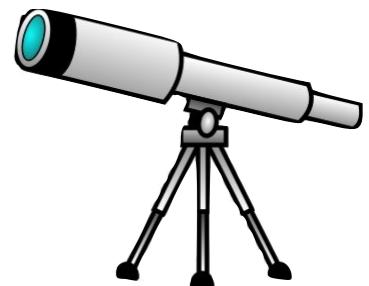


# Neural Network

# Renormalization Group



Lei Wang (王磊)

Institute of Physics, CAS

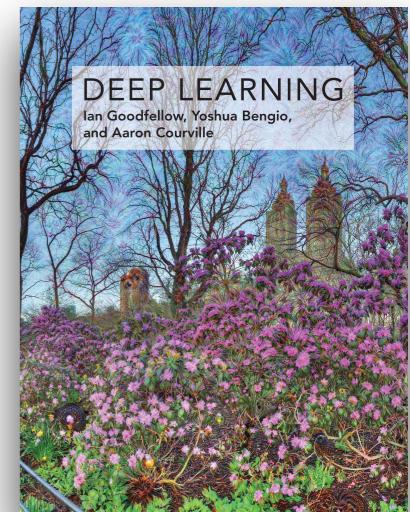
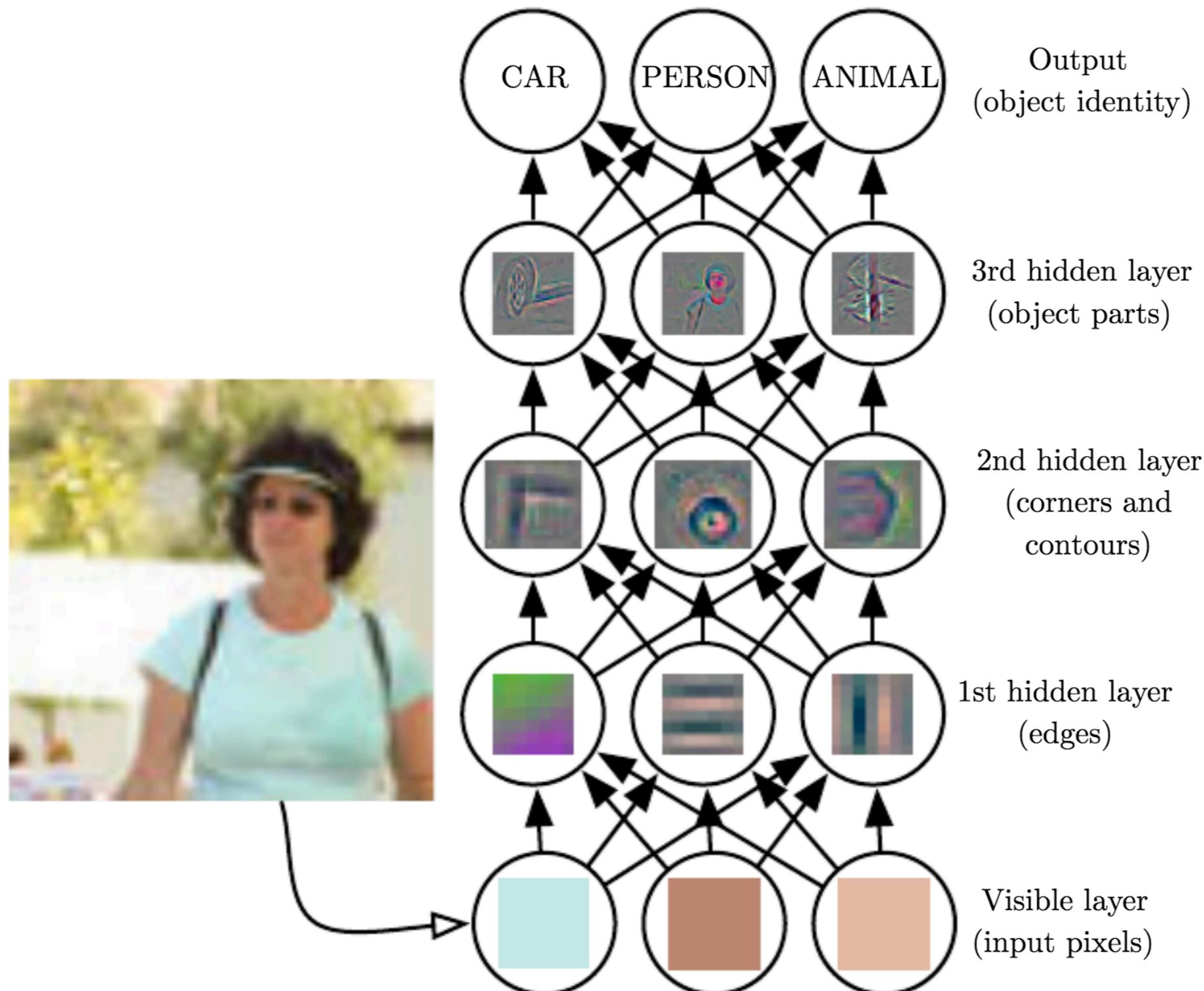
<https://wangleiphy.github.io>

# Collaborator



Shuo-Hui Li (李砾辉)  
PhD. student at IOP, CAS

# Deep Neural Network and RG



# Deep learning and the renormalization group

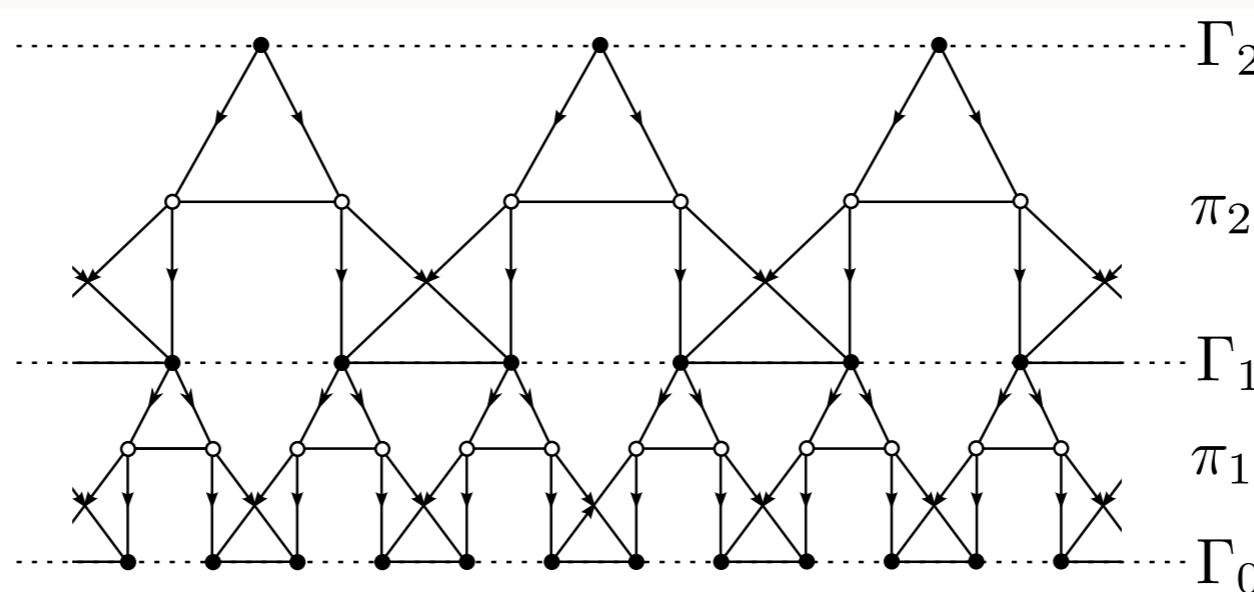


Cédric Bény

15 Jan 2013 ICLR 2013 conference submission readers: everyone

**Decision:** reject

**Abstract:** Renormalization group methods, which analyze the way in which the effective behavior of a system depends on the scale at which it is observed, are key to modern condensed-matter theory and particle physics. The aim of this paper is to compare and contrast the ideas behind the renormalization group (RG) on the one hand and deep machine learning on the other, where depth and scale play a similar role. In order to illustrate this connection, we review a recent numerical method based on the RG---the multiscale entanglement renormalization ansatz (MERA)---and show how it can be converted into a learning algorithm based on a generative hierarchical Bayesian network model. Under the assumption---common in physics---that the distribution to be learned is fully characterized by local correlations, this algorithm involves only explicit evaluation of probabilities, hence doing away with sampling.



arxiv:1301.3124

# Deep learning and the renormalization group



Cédric Bény

15 Jan 2013 ICLR 2013 conference submission readers: everyone

**Decision:** reject

Yann LeCun

05 Apr 2013 ICLR 2013 submission review readers: everyone

**Review:** It seems to me like there could be an interesting connection between approximate inference in graphical models and the renormalization methods.

There is in fact a long history of interactions between condensed matter physics and graphical models. For example, it is well known that the loopy belief propagation algorithm for inference minimizes the Bethe free energy (an approximation of the free energy in which only pairwise interactions are taken into account and high-order interactions are ignored). More generally, variational methods inspired by statistical physics have been a very popular topic in graphical model inference.

The renormalization methods could be relevant to deep architectures in the sense that the grouping of random variable resulting from a change of scale could be made analogous with the pooling and subsampling operations often used in deep models.

It's an interesting idea, but it will probably take more work (and more tutorial expositions of RG) to catch the attention of this community.

# A Common Logic to Seeing Cats and Cosmos

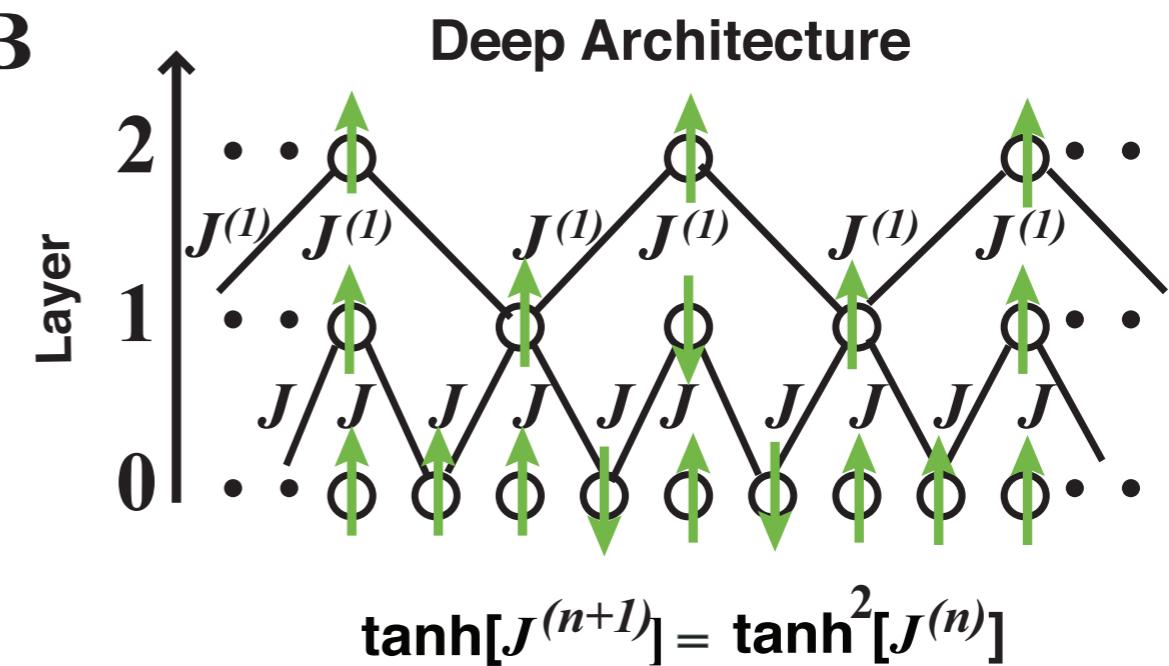
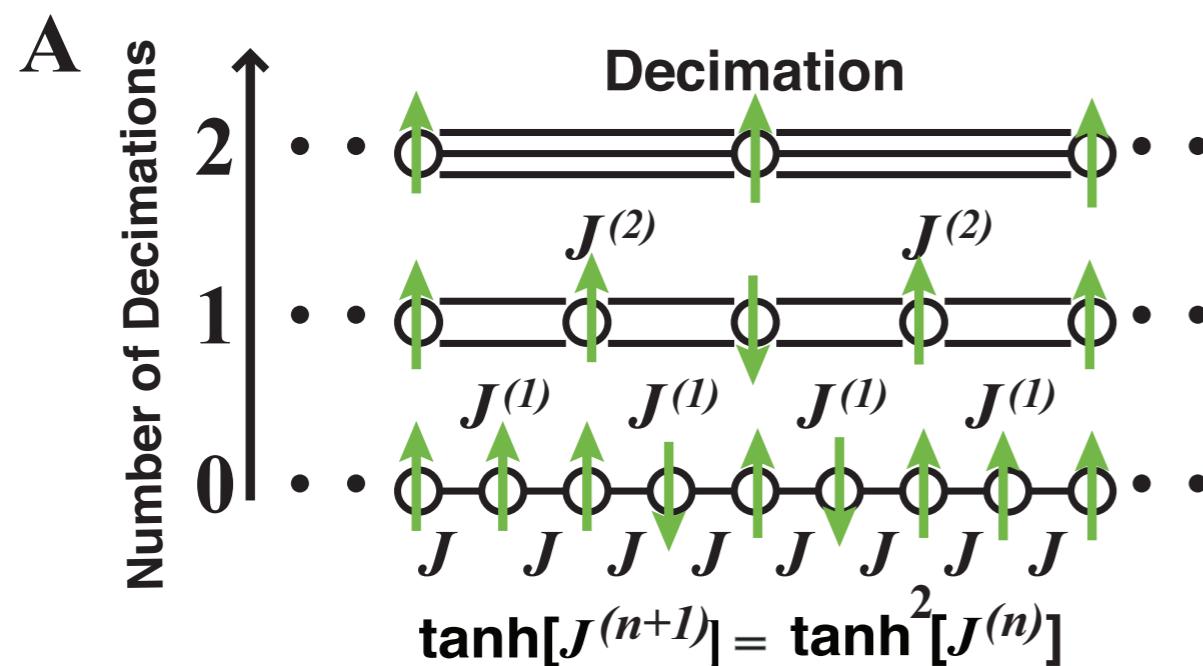


Olena Shmahalo / Quanta Magazine

There may be a universal logic to how physicists, computers and brains tease out important features from among other irrelevant bits of data.

“An exact mapping between the Variational Renormalization Group and Deep Learning”, Mehta and Schwab, 1410.3831

# “Exact Mapping”



$$e^{-H(\mathbf{h})} = \sum_{\mathbf{x}} e^{T(\mathbf{x}, \mathbf{h}) - H(\mathbf{x})}$$

RG Transformation

$$e^{-E(\mathbf{h})} = \sum_{\mathbf{x}} e^{-E(\mathbf{x}, \mathbf{h})}$$

Boltzmann Machine

# Harsh comments below the Quanta Magazine article

Noah says:

December 26, 2014 at 9:54 am

I just spend an hour reading Mehta-Schwab paper from the beginning to end. Let me say that “A Common Logic to Seeing Cats and Cosmos” is a sensationalist article about a trivial paper, which will have no impact whatsoever. The whole M-S paper is based on the fact that couplings of two systems appear in more than one context and that distributions can sometimes appear as marginal distributions on product spaces. There is no one-to-one mappings between renormalization group (RG) scheme of Kadanoff and Restricted Boltzmann Machines (RBM) in Deep Neural Networks (DNN) in their paper. What they show is that RBM can be represented as a RG scheme with a very specific choice of coupling function  $T$  in equation (18). Conveniently, this coupling function depends on the Hamiltonian of the spin system, which it normally should not. Equivalence in equations (8) and (9) is also not correct. Condition (9) of course implies that the scheme is exact, but not the other way around, unless the authors make some implicit assumptions about coupling function  $T$  not mentioned in the paper. The paper contains no non-trivial ideas, it does not “open up a door to something very exciting”, and I will not hold my breath expecting new breakthroughs because of this connection.

# Dictionary: RG vs Deep Learning

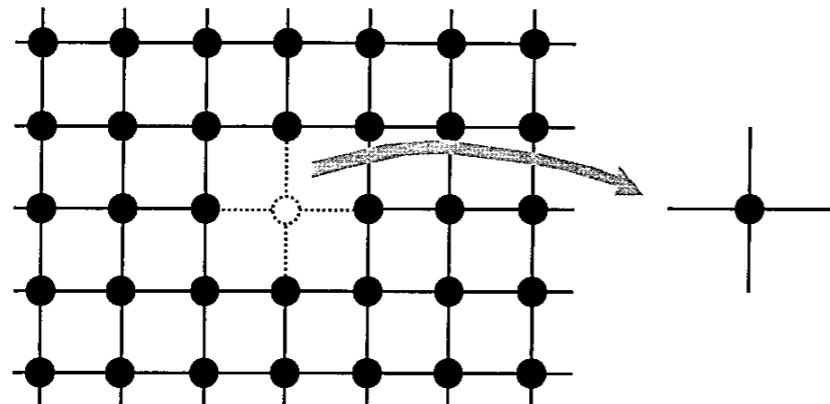
Property	Variational RG	Deep Belief Networks
How input distribution is defined	Hamiltonian defining $P(v)$	Data samples drawn from $P(v)$
How interactions are defined	$T(v,h)$	$E(v,h)$
Exact transformation	$\text{Tr}_h e^{T(v,h)} = 1$	KL divergence between $P(v)$ and variational distribution is zero
Approximations	Minimize or bound free energy differences	Minimize the KL divergence
Method	Analytic (mostly)	Numerical
What happens under coarse-graining	Relevant operators grow/irrelevant shrink	New features emerge

# More on the DL-RG Connections

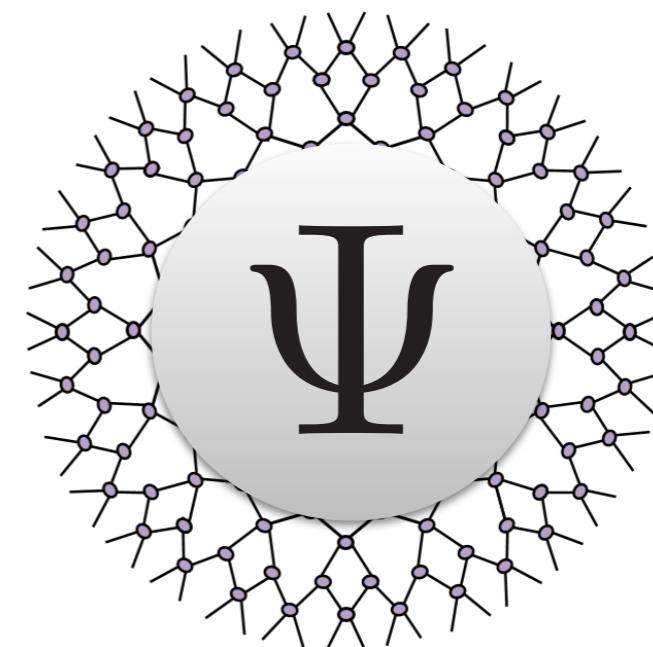
- “Why does deep and cheap learning work so well ”,  
Lin, Tegmark, Rolnick, 1608.08225
- Comment on the above paper, Schwab and Mehta,  
1609.03541
- PCA meets RG, Bradde and Bialek, 1610.09733
- Mutual information RG, Koch-Janusz and Ringel,  
1704.06279
- Media coverage/blog posts/student term papers etc

# Physics-DL connection is more general than we thought

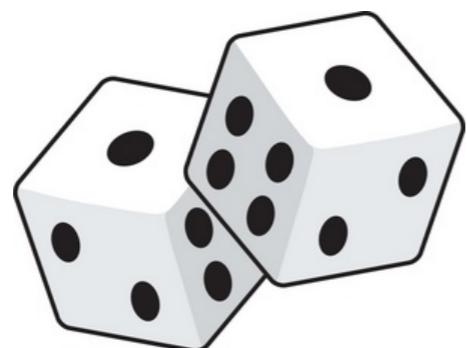
## Variational Mean Field



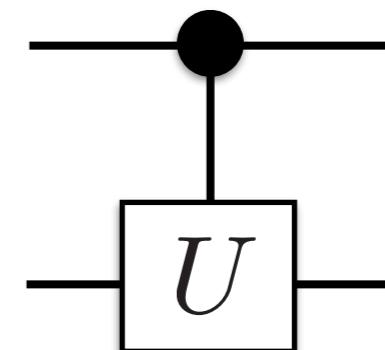
## Tensor Networks



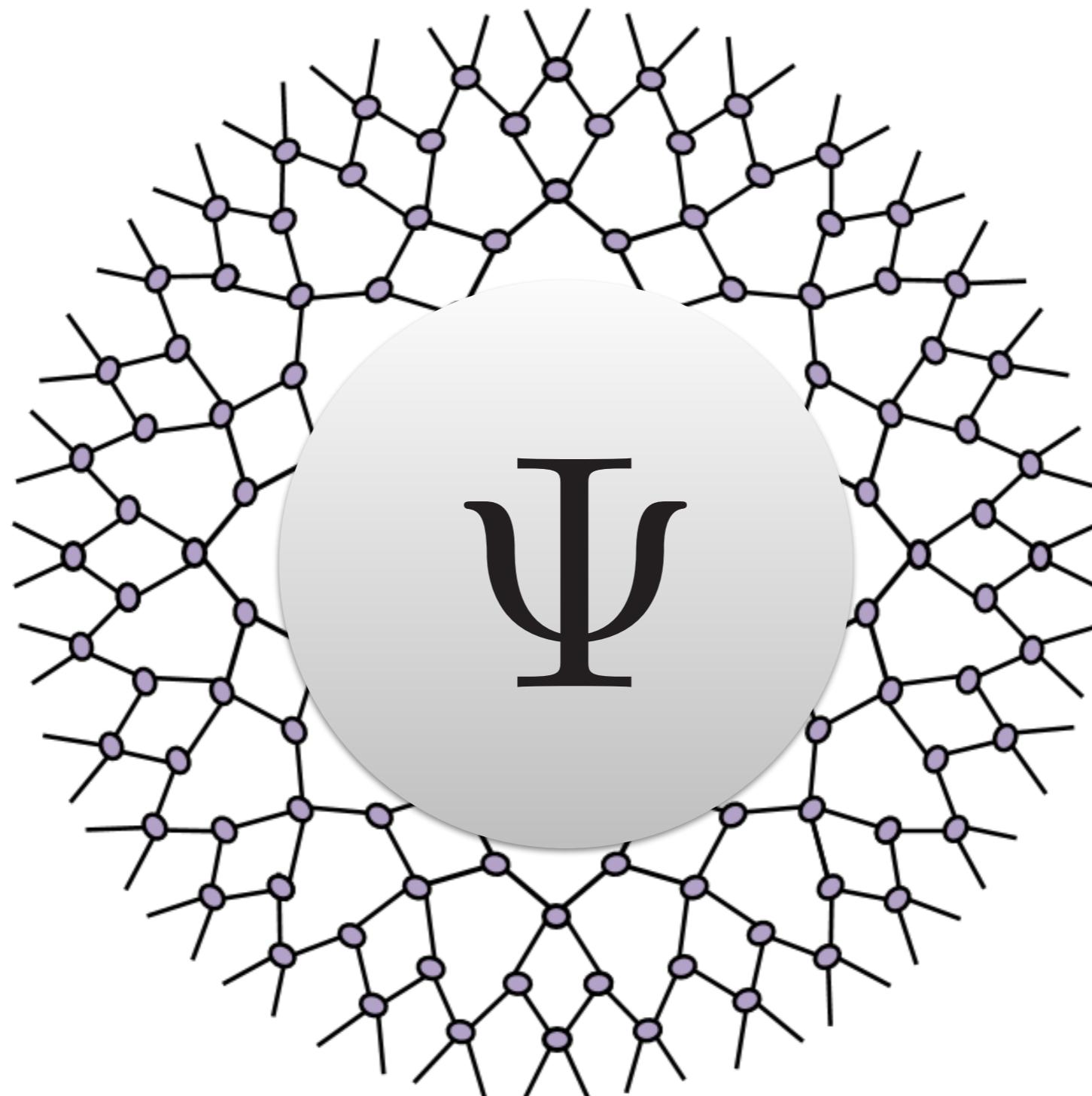
## Monte Carlo Methods



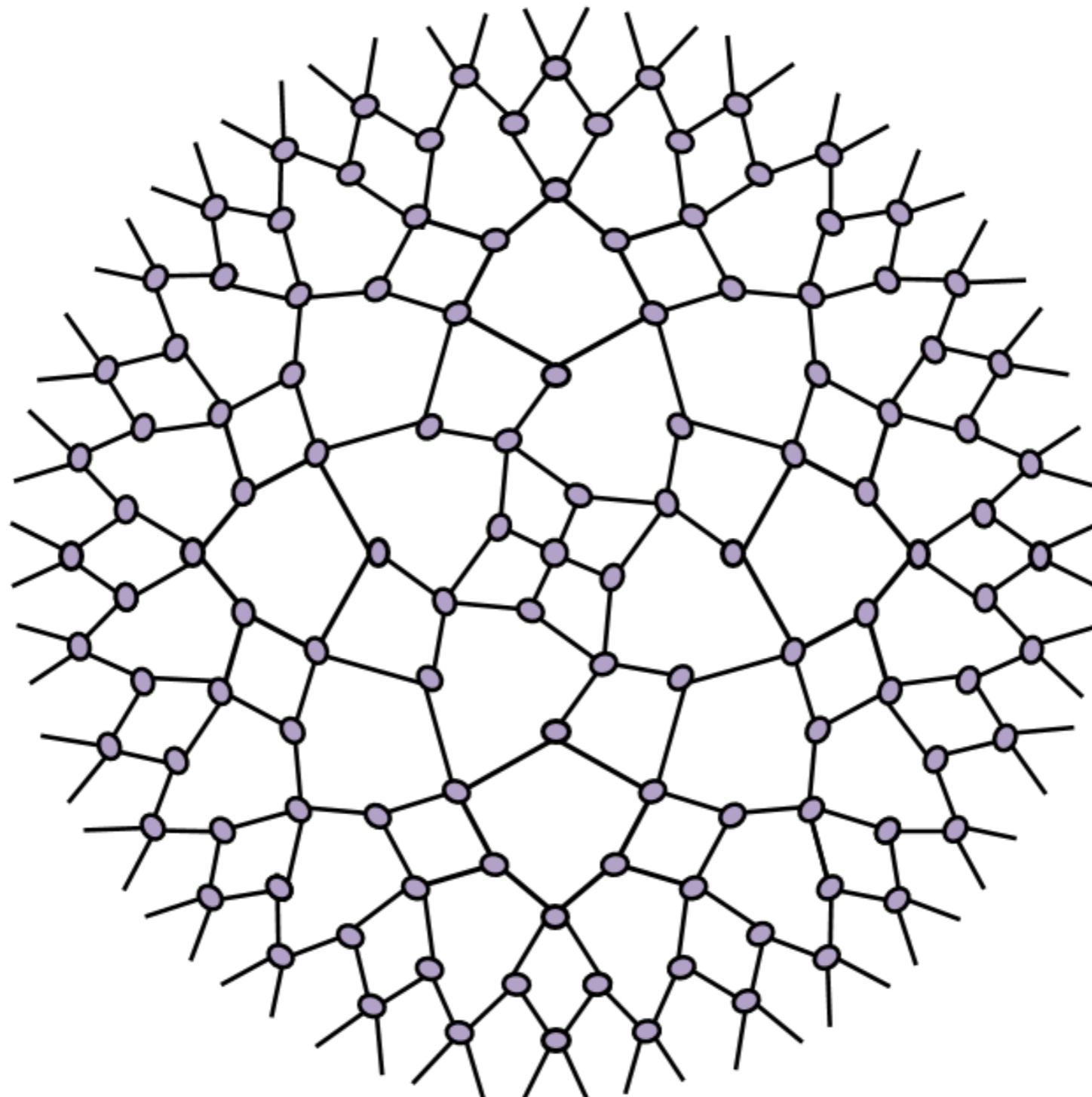
## Quantum Computing



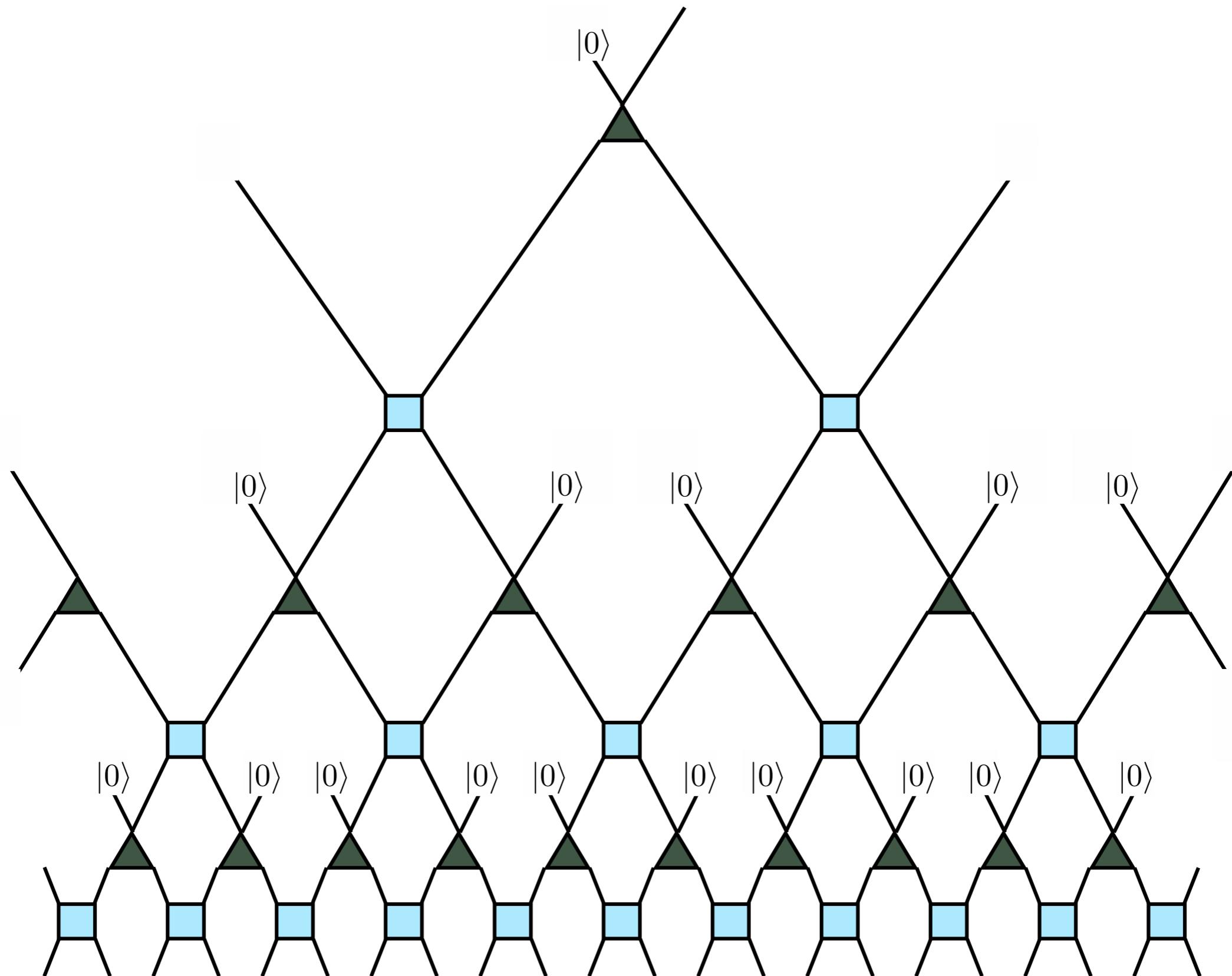
# Multi-Scale Entanglement Renormalization Ansatz



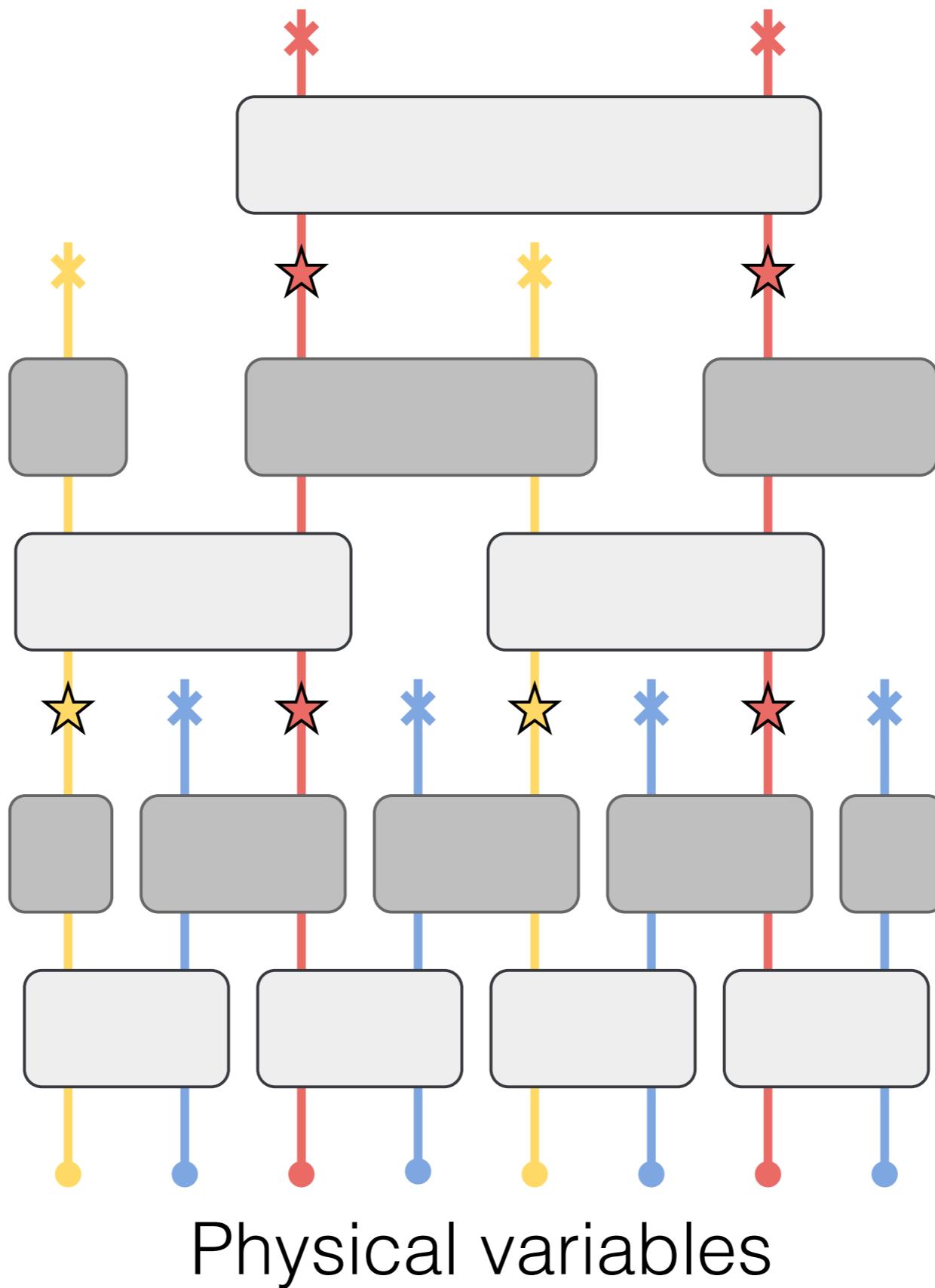
# Multi-Scale Entanglement Renormalization Ansatz



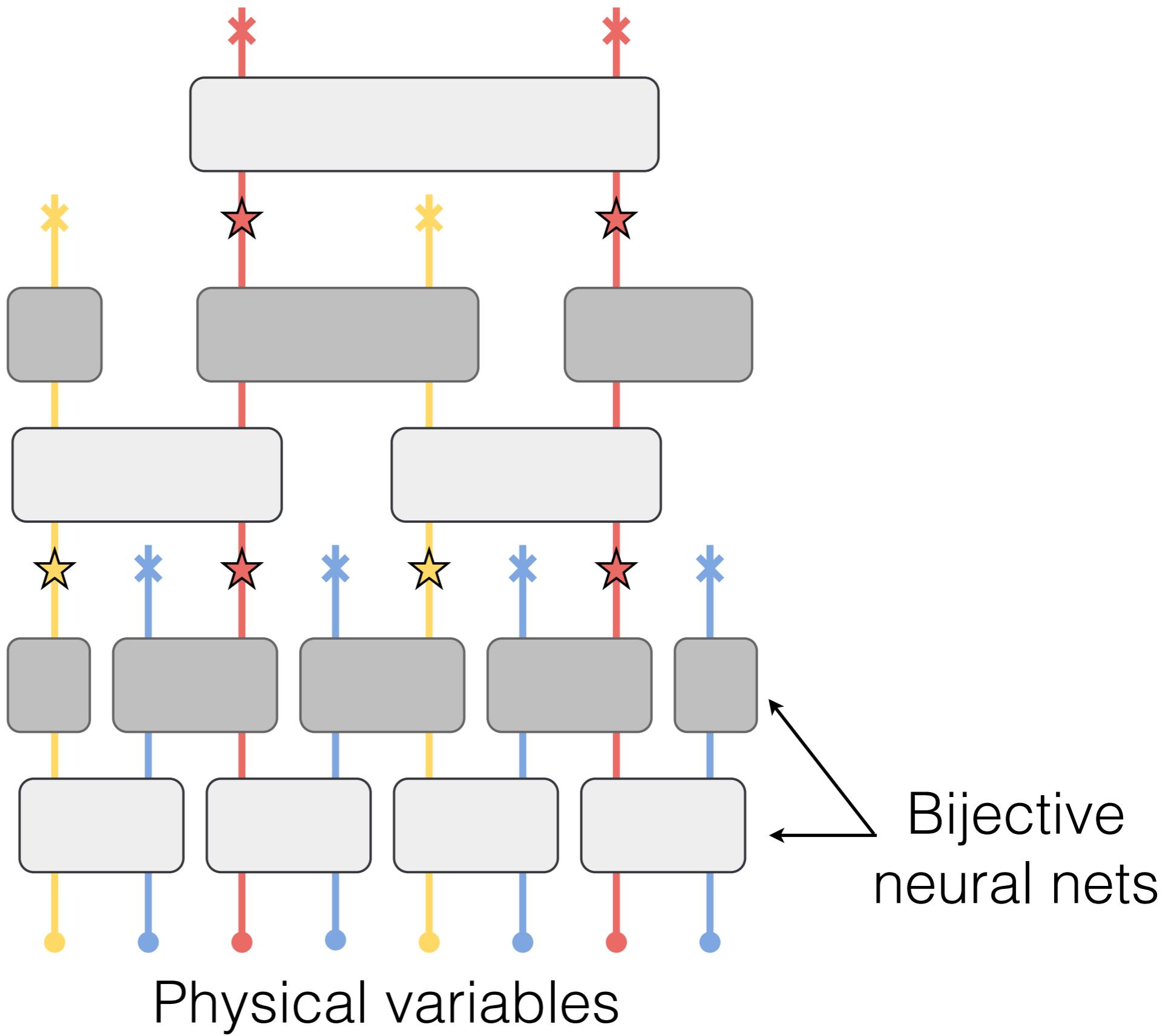
# MERA as a quantum circuit



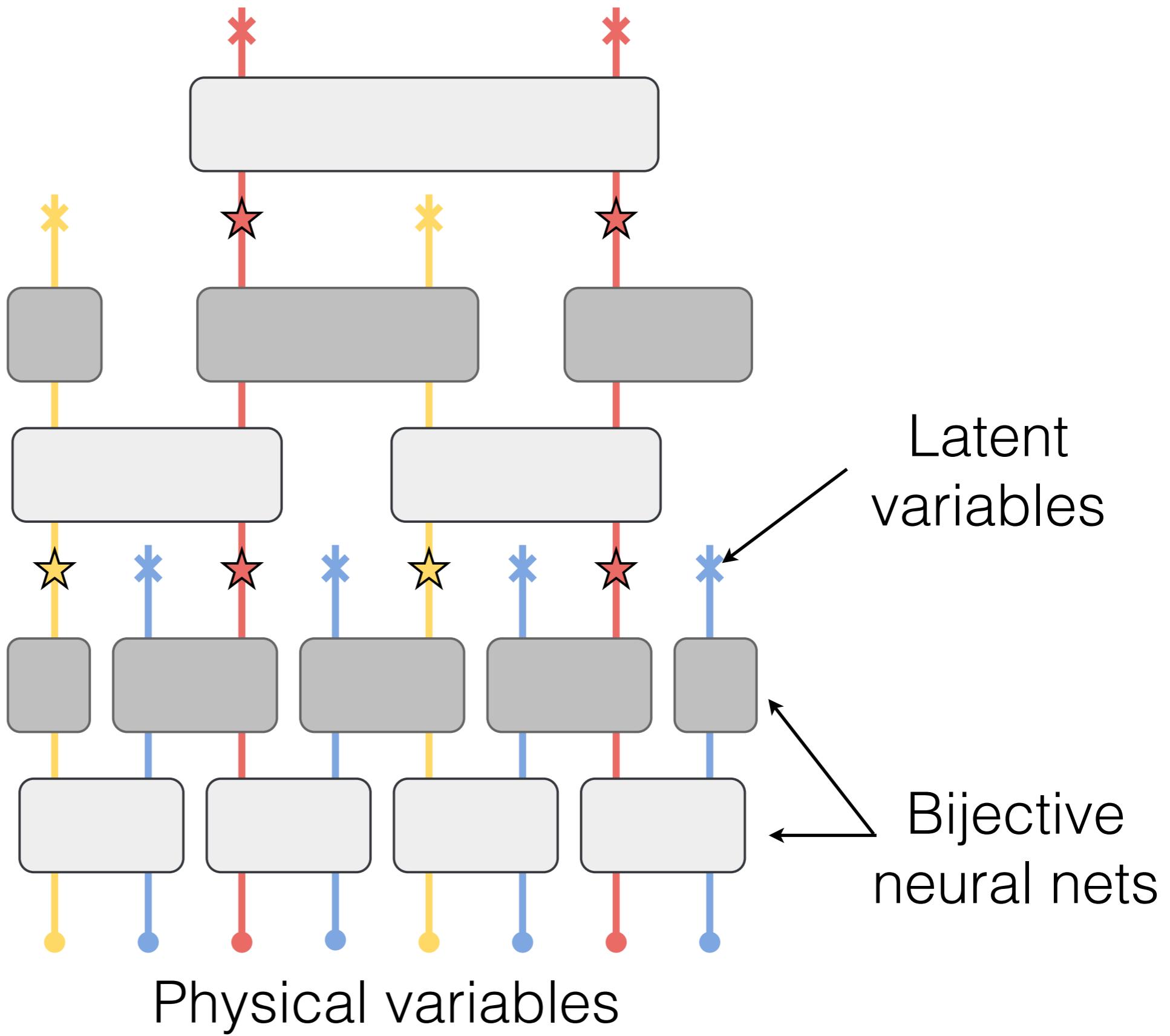
# Neural Network Renormalization Group



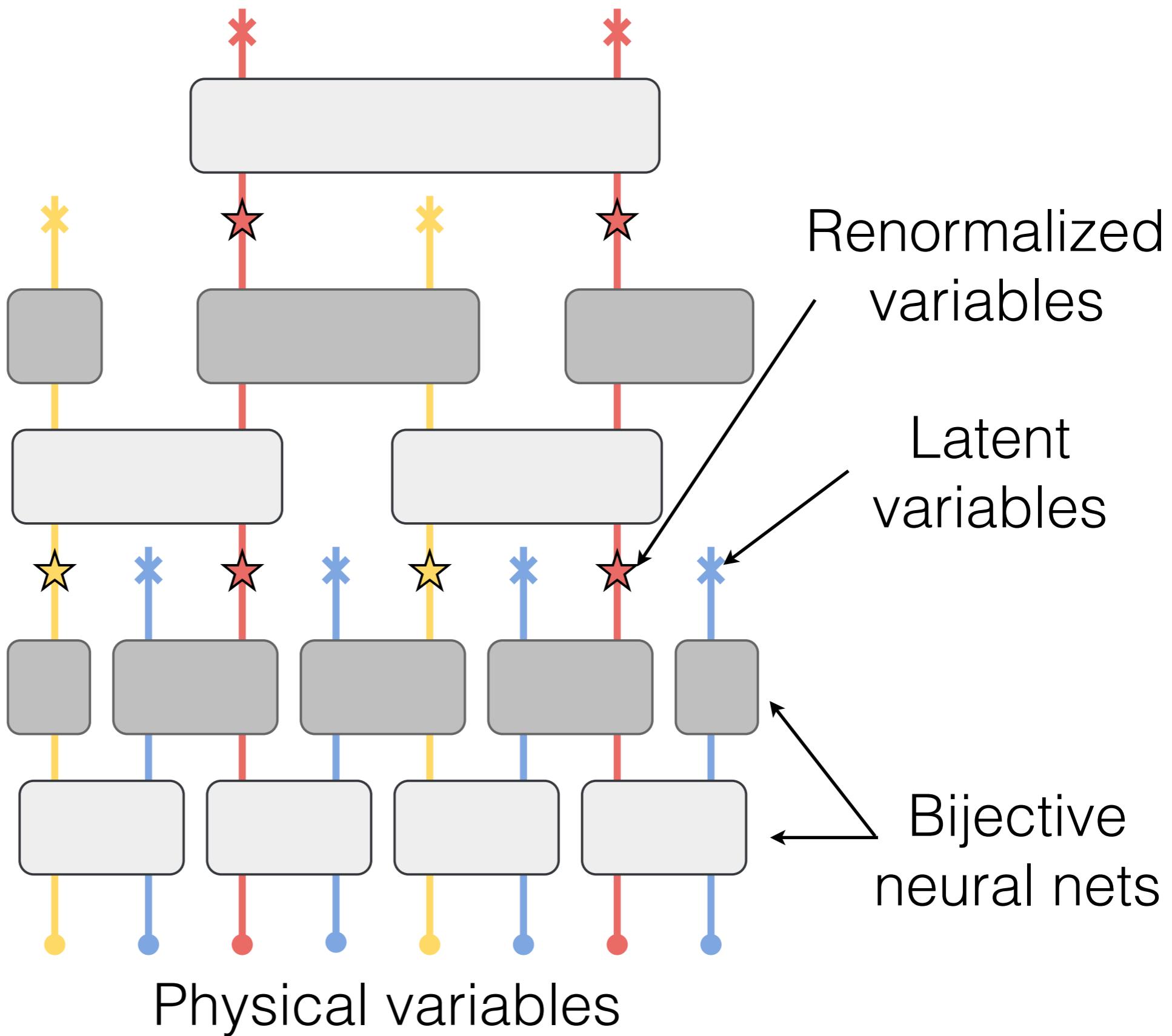
# Neural Network Renormalization Group



# Neural Network Renormalization Group



# Neural Network Renormalization Group

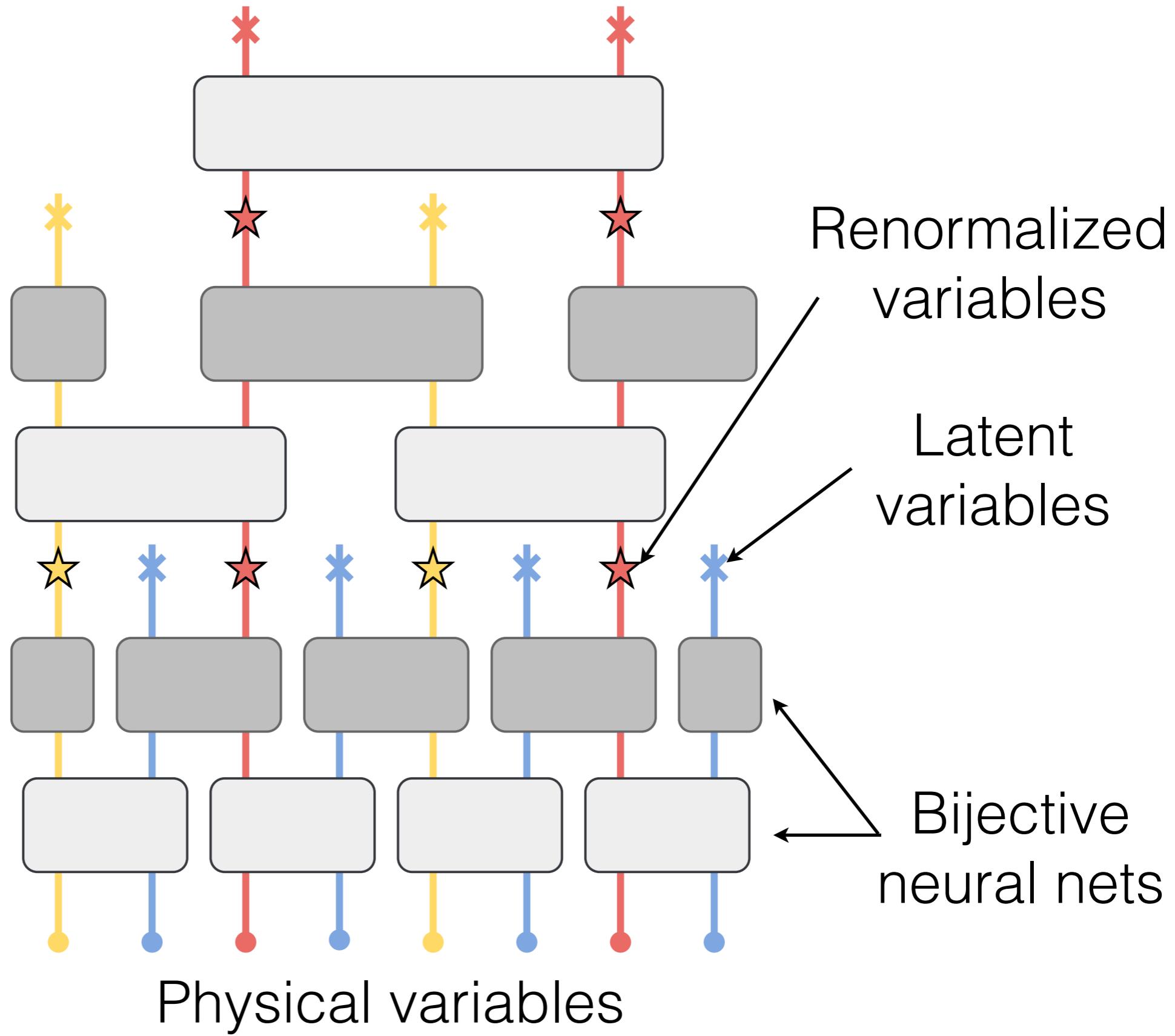


# Neural Network Renormalization Group

$$z = g^{-1}(x)$$

Inference  
Generate

$$x = g(z)$$

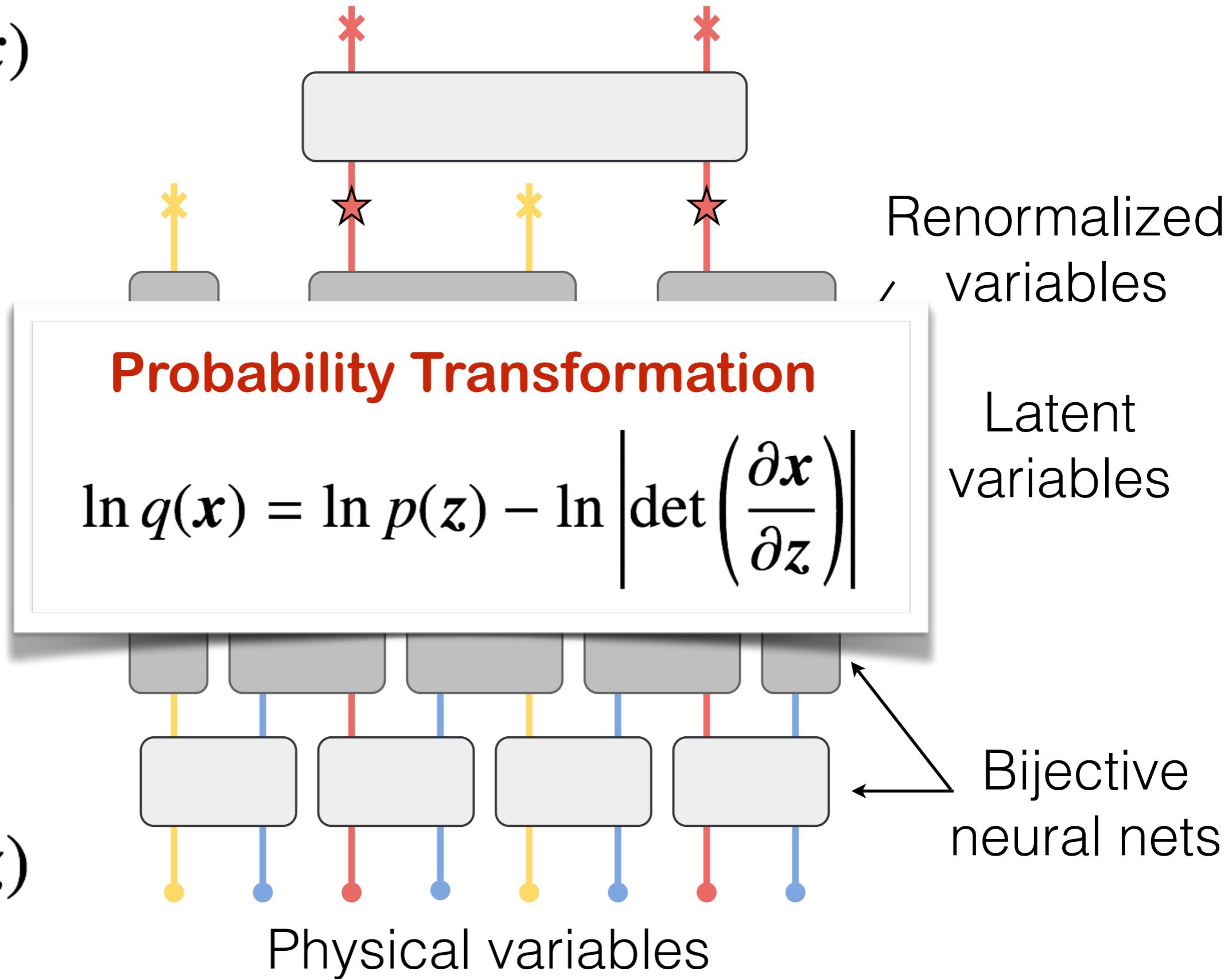


# Neural Network Renormalization Group

$$z = g^{-1}(x)$$

Inference  
Generate

$$x = g(z)$$



# Bijector Block

Bijective & Differentiable map, i.e., [Diffeomorphism](#)

Forward

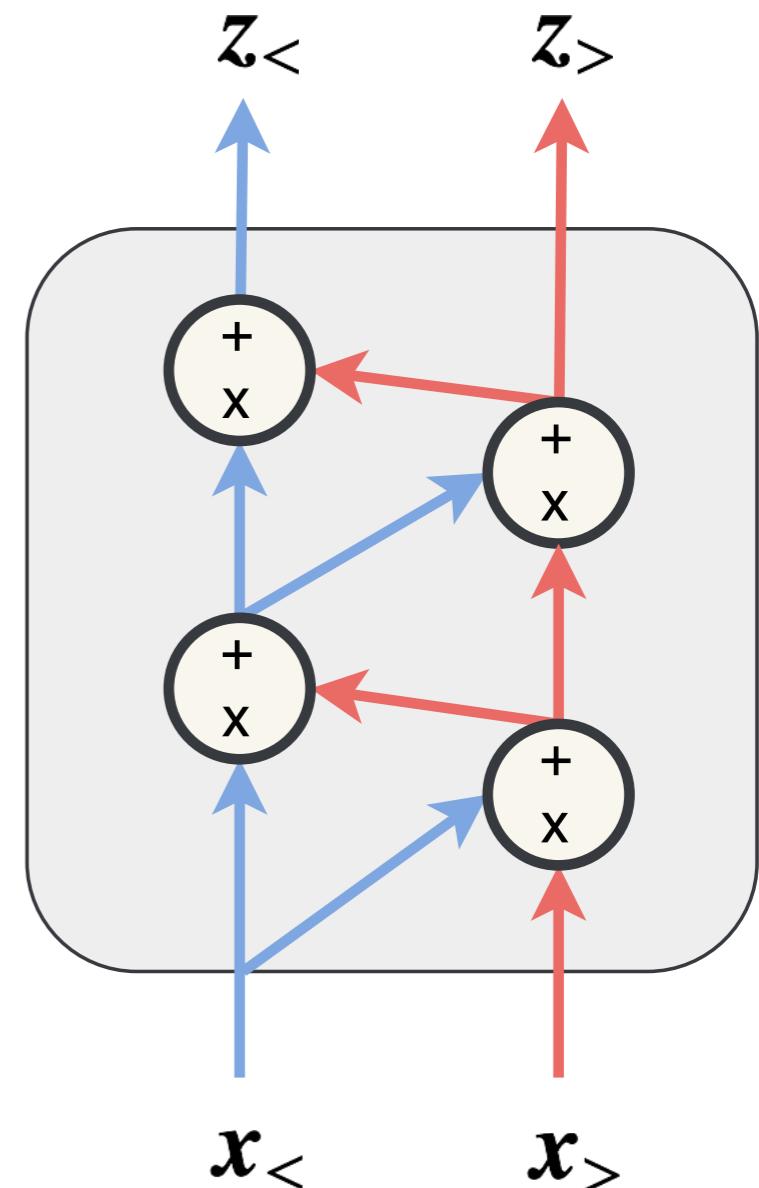
$$\begin{cases} \mathbf{x}_< = \mathbf{z}_< \\ \mathbf{x}_> = \mathbf{z}_> \odot e^{s(\mathbf{z}_<)} + t(\mathbf{z}_<) \end{cases}$$

Backward

$$\begin{cases} \mathbf{z}_< = \mathbf{x}_< \\ \mathbf{z}_> = (\mathbf{x}_> - t(\mathbf{x}_<)) \odot e^{-s(\mathbf{x}_<)} \end{cases}$$

Log-Det-Jacobian

$$\ln \left| \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i [s(\mathbf{z}_<)]_i$$



Normalizing flow, Rezende et al, 1505.05770  
Special case: Real NVP, Dinh et al, 1605.08803

# Bijector Block

Bijective & Differentiable map, i.e., [Diffeomorphism](#)

Forward

Arbitrary  
neural nets

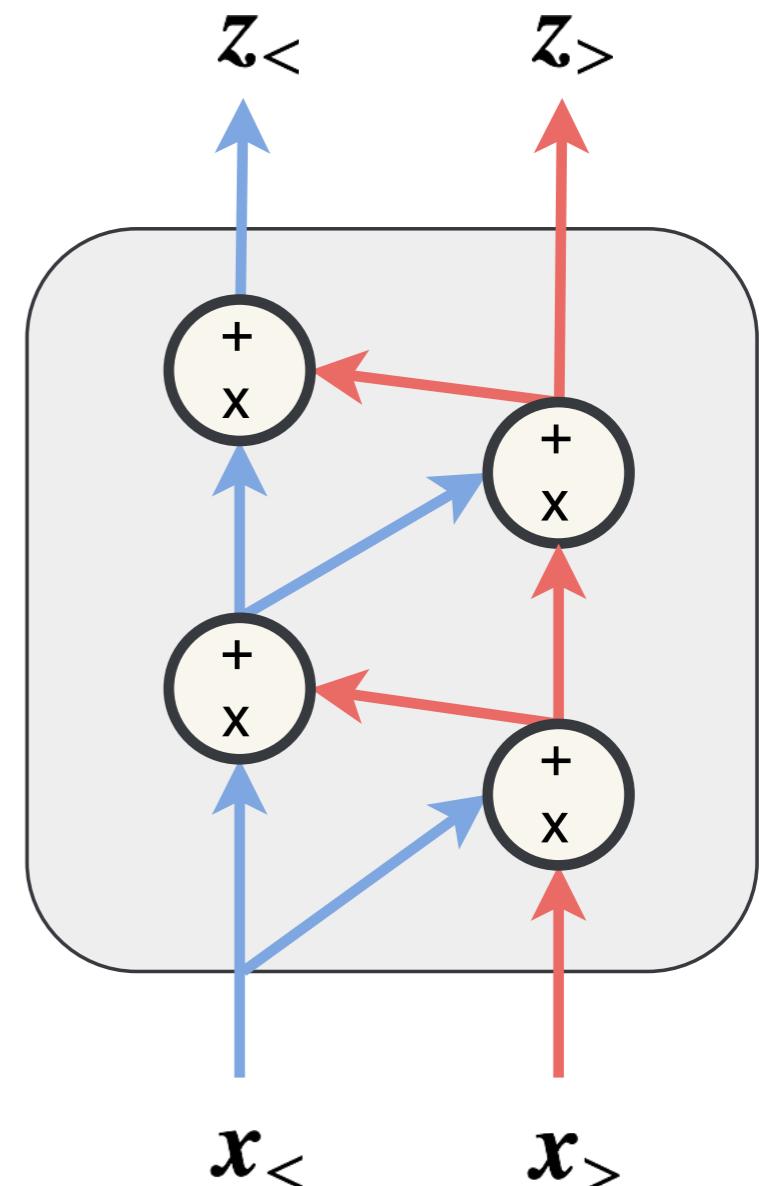
$$\begin{cases} \mathbf{x}_< = \mathbf{z}_< \\ \mathbf{x}_> = \mathbf{z}_> \odot e^{s(\mathbf{z}_<)} + t(\mathbf{z}_<) \end{cases}$$

Backward

$$\begin{cases} \mathbf{z}_< = \mathbf{x}_< \\ \mathbf{z}_> = (\mathbf{x}_> - t(\mathbf{x}_<)) \odot e^{-s(\mathbf{x}_<)} \end{cases}$$

Log-Det-Jacobian

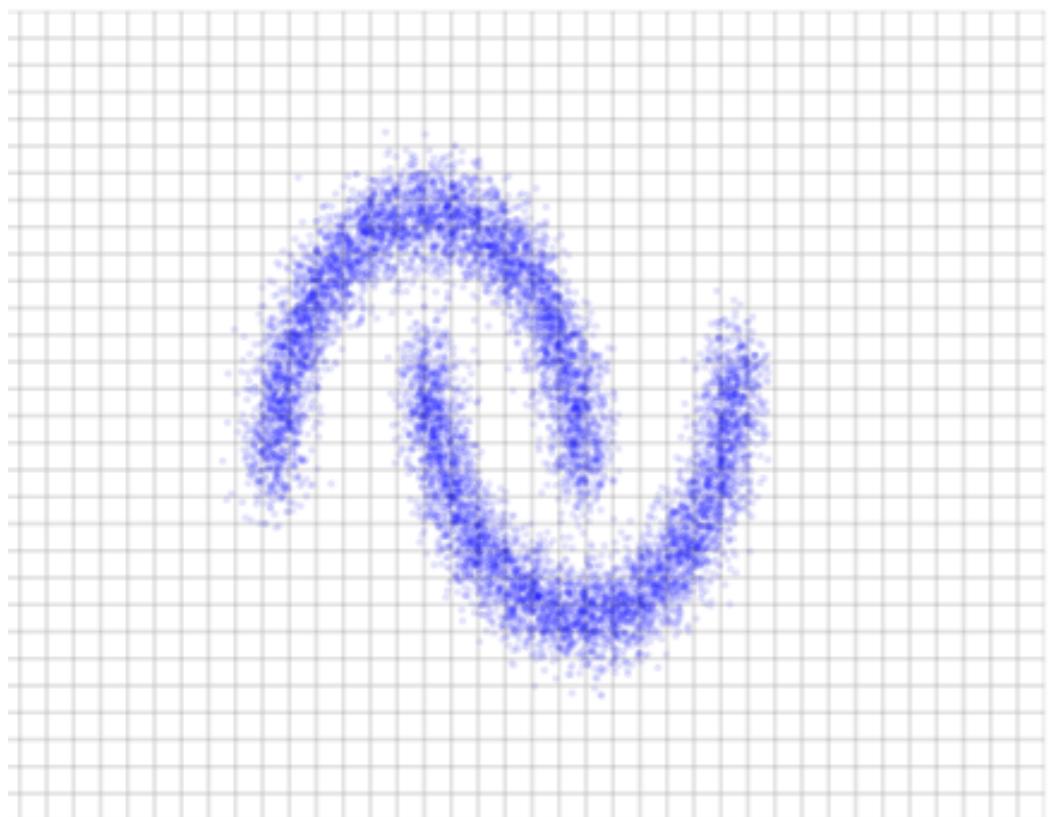
$$\ln \left| \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i [s(\mathbf{z}_<)]_i$$



Normalizing flow, Rezende et al, 1505.05770  
Special case: Real NVP, Dinh et al, 1605.08803

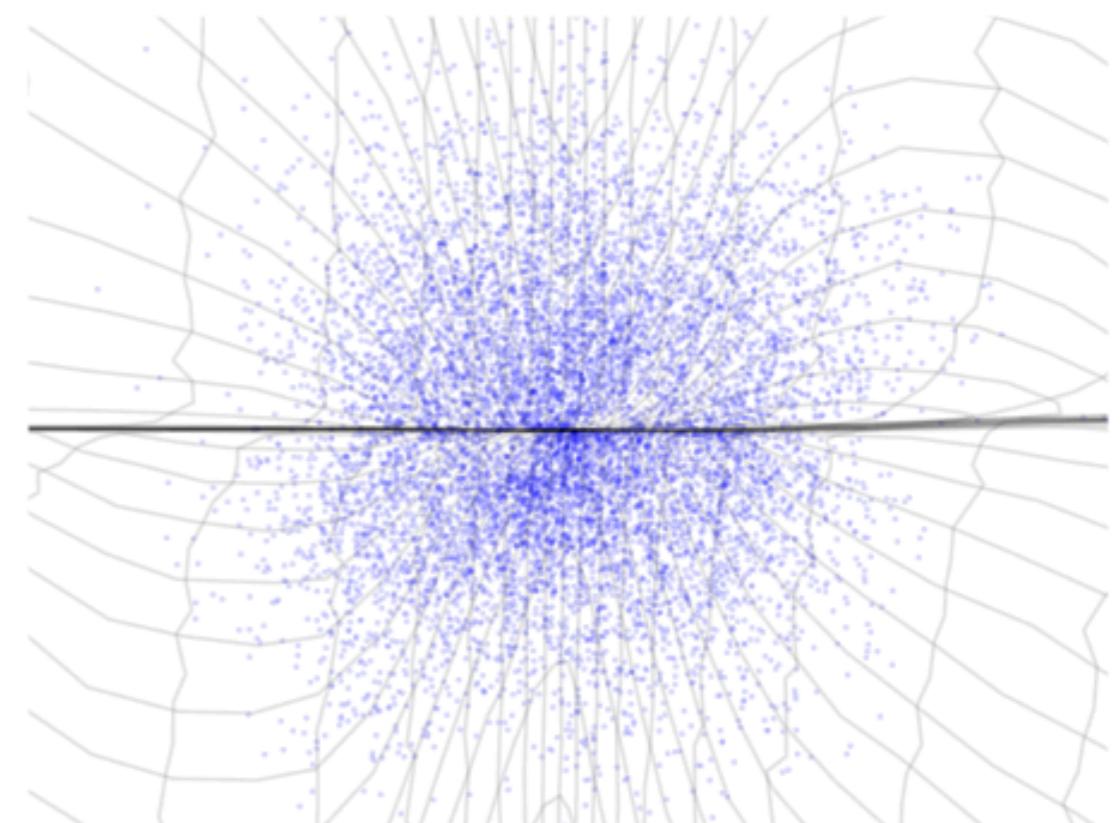
# Inference

Data space



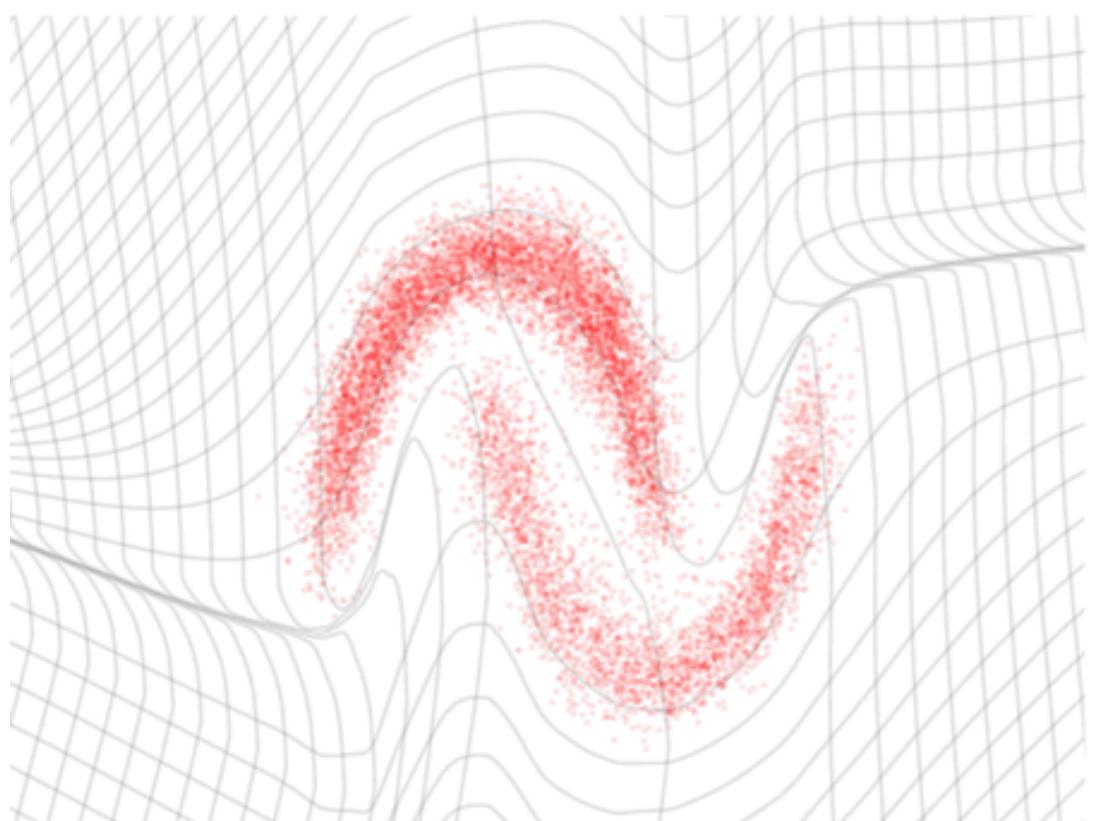
$$g^{-1}(\mathbf{x}) = \mathbf{z}$$

Latent space



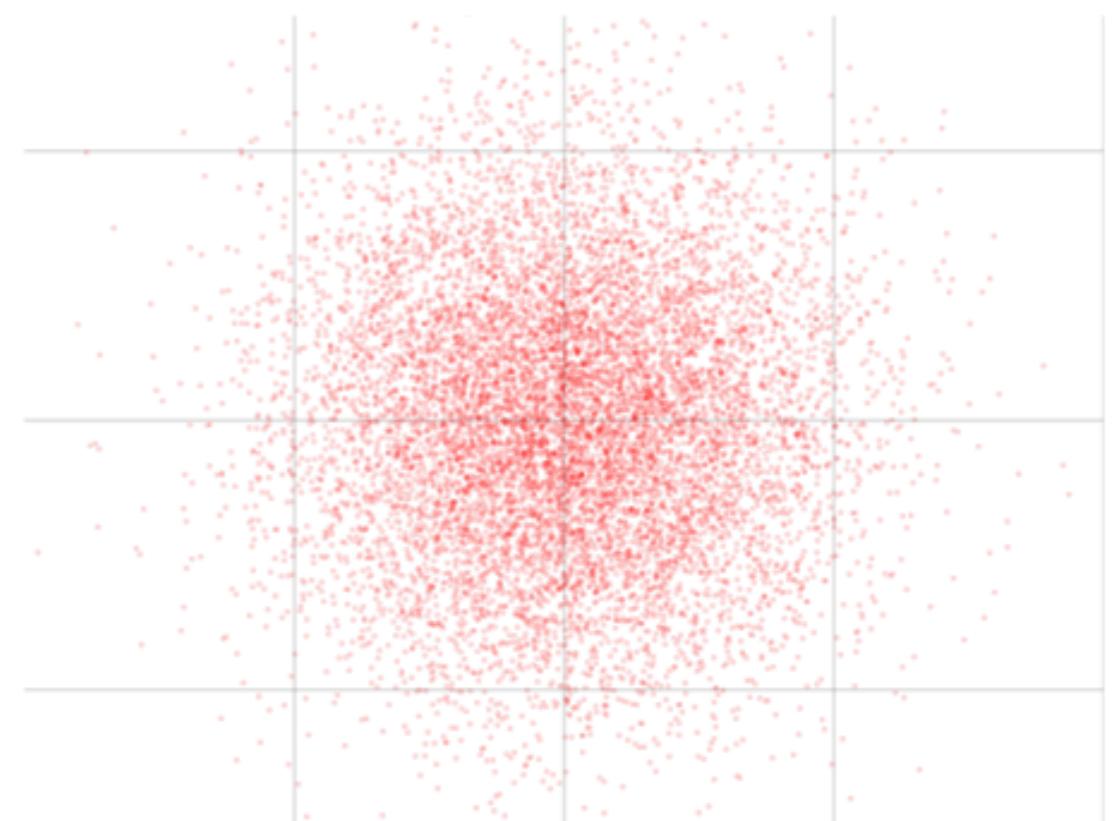
# Generate

Data space



$$x = g(z)$$

Latent space

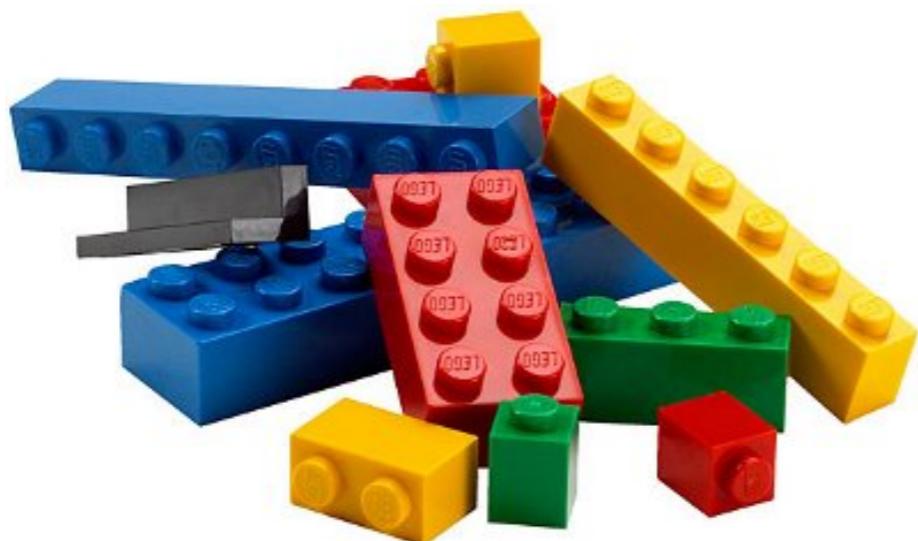


# Bijectors form a group

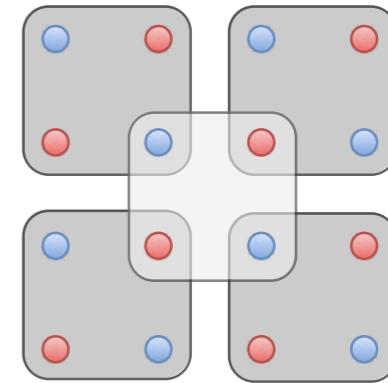
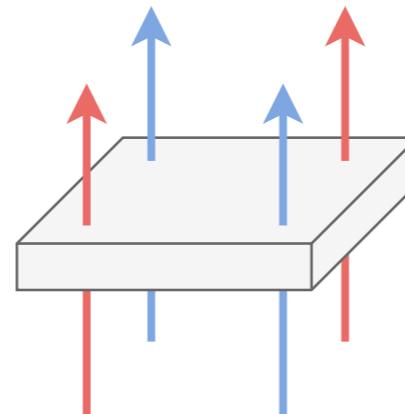
$$\mathbf{x} = g(\mathbf{z})$$

$$g = \cdots \circ g_2 \circ g_1$$

$$\ln \left| \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{z}} \right) \right| = \sum_i \ln \left| \det \left( \frac{\partial g_{i+1}}{\partial g_i} \right) \right|$$

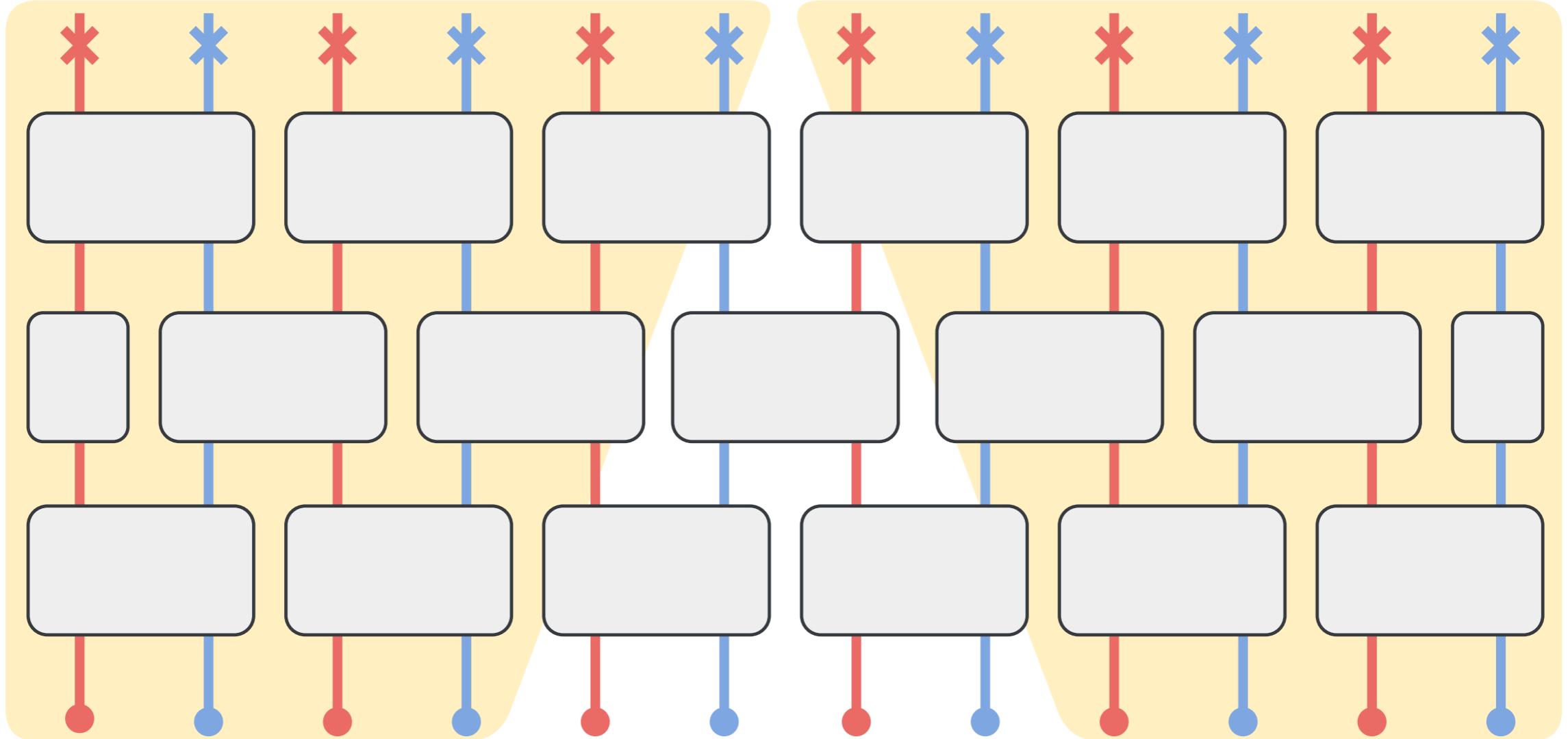


Modular design



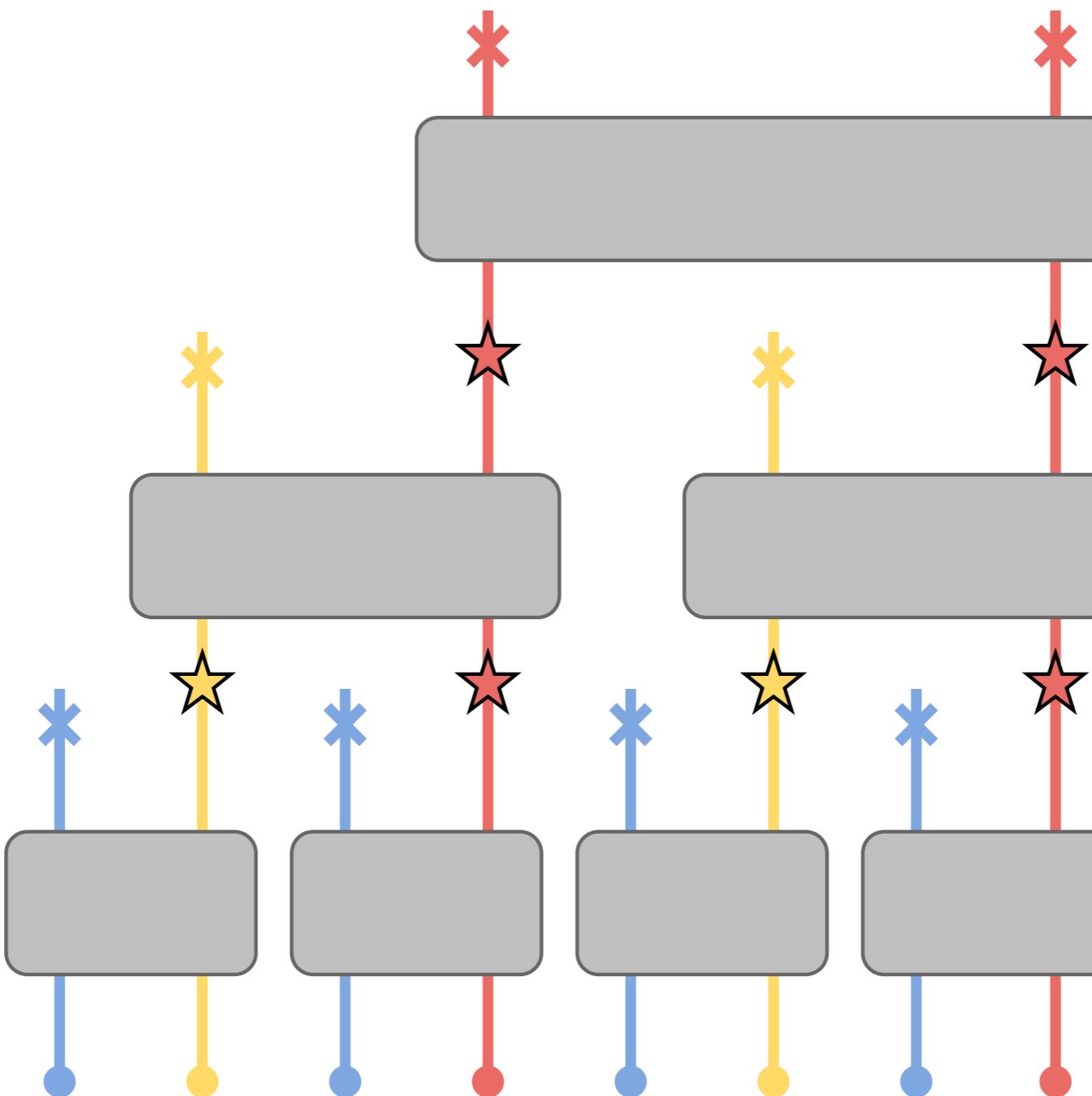
Flexible blocks and stacking

# “Disentangler” only architecture

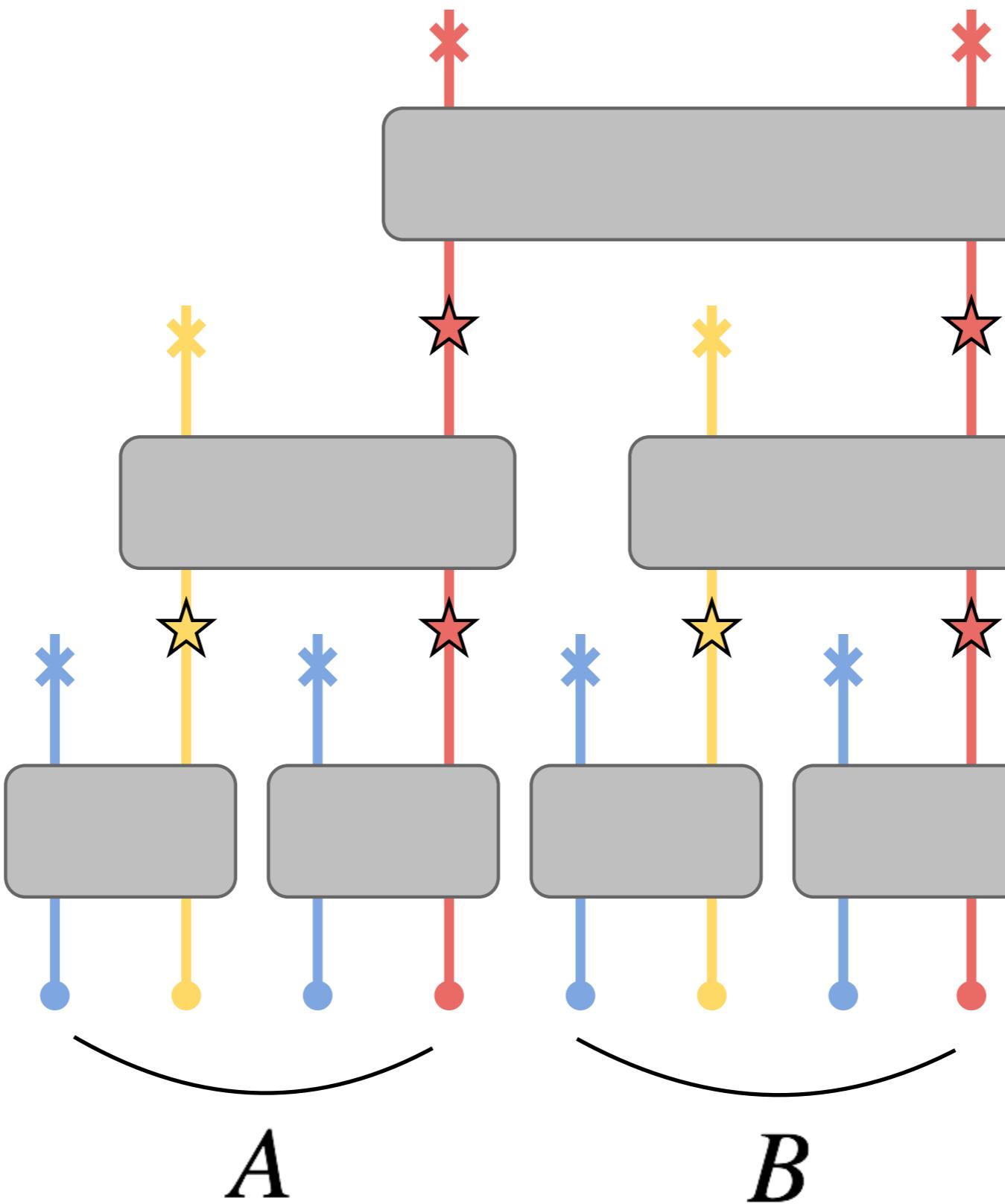


Correlation length  $\sim$  Network depth

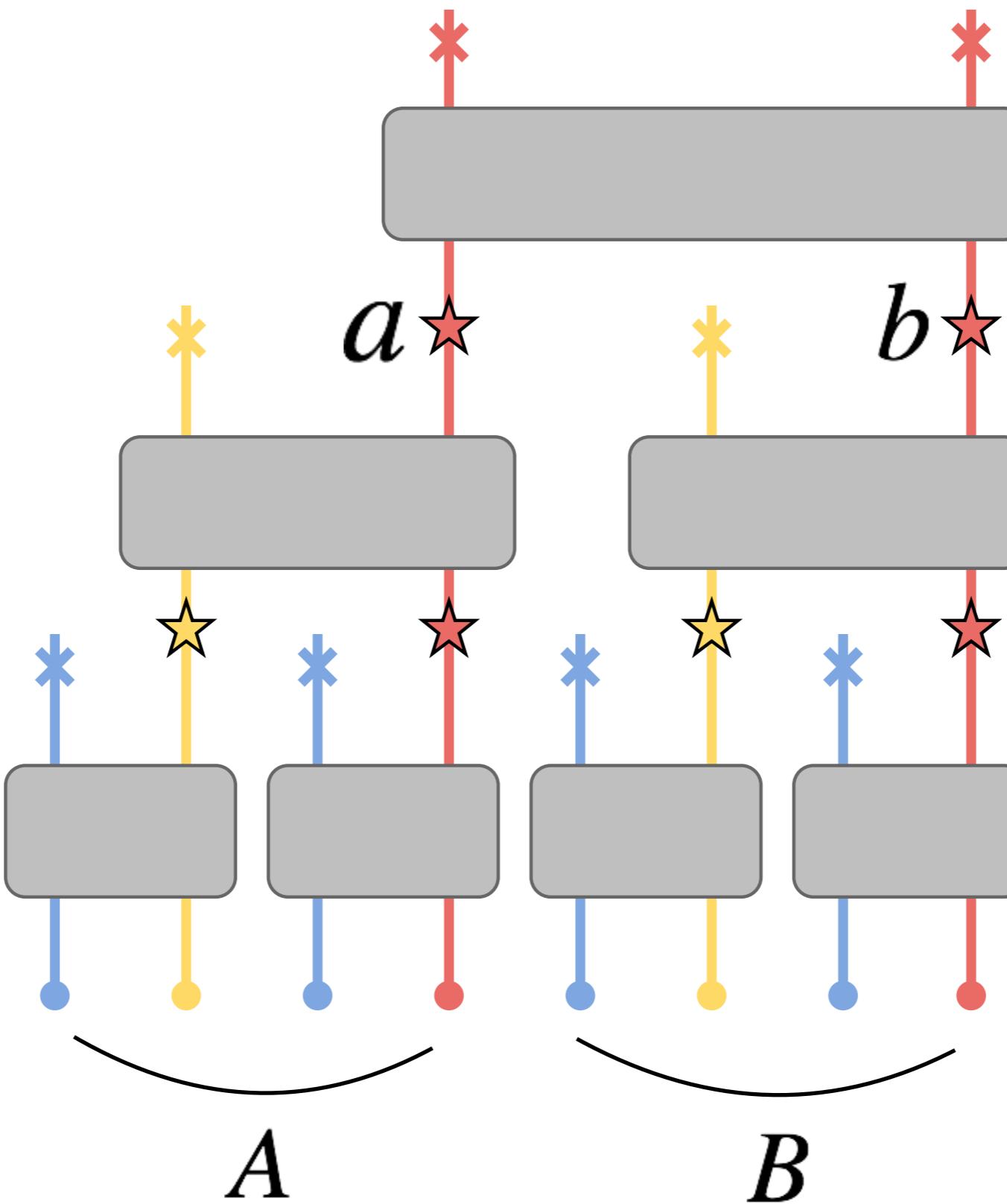
# “Decimator” only architecture



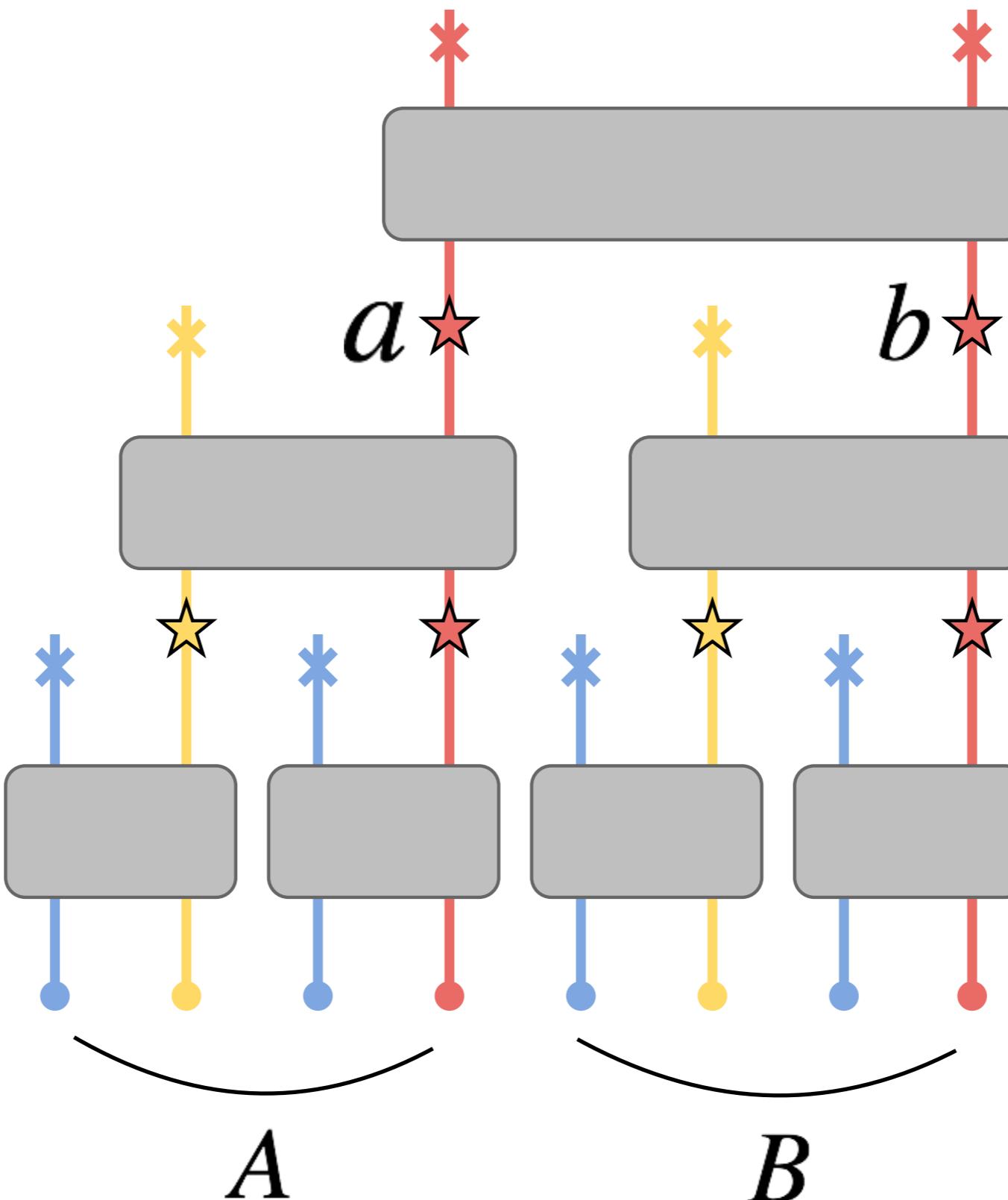
# “Decimator” only architecture



# “Decimator” only architecture



# “Decimator” only architecture



$$I(A : B) = I(a : b)$$

Mutual Information Bottleneck

# Probability Density *Estimation*

Given a dataset, learn its probability density by minimizing the negative likelihood

$$\text{NLL}_{\theta} = - \sum_{x \in \text{dataset}} \ln q_{\theta}(x)$$

Network parameters

Equivalent to reduce the Kullback–Leibler divergence

$$\text{KL}\left(\frac{\pi(x)}{Z} \parallel q_{\theta}(x)\right)$$

However, typical Stat-Mech problem has access only to the energy function

# Probability Density *Distillation*

Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$

# Probability Density *Distillation*

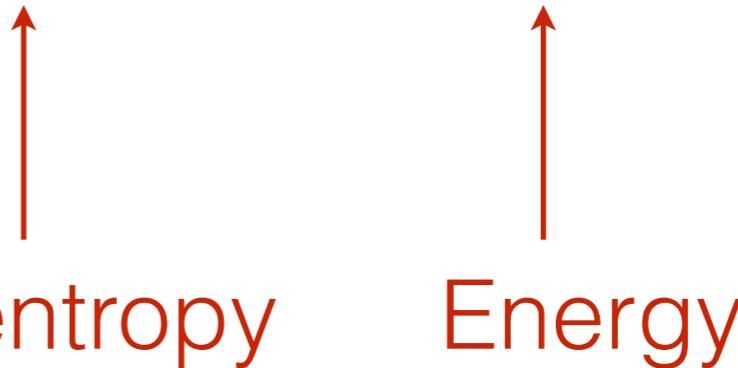
Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$

↑  
Negentropy

# Probability Density *Distillation*

Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$


Negentropy      Energy

# Probability Density *Distillation*

Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$

↑                      ↑                      ↑

“Variational  
Free Energy”              Negentropy              Energy

# Probability Density *Distillation*

Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$

“Variational  
Free Energy”

Negentropy

Energy

$$Z = \int d\mathbf{x} \pi(\mathbf{x})$$

Partition function

# Probability Density *Distillation*

Learn from the data generated by the network itself

$$\mathcal{L}_\theta = \int d\mathbf{x} q_\theta(\mathbf{x}) [\ln q_\theta(\mathbf{x}) - \ln \pi(\mathbf{x})]$$

↑                      ↑                      ↑

“Variational  
Free Energy”              Negentropy              Energy

$$\mathcal{L}_\theta + \ln Z = \text{KL}\left(q_\theta(\mathbf{x}) \parallel \frac{\pi(\mathbf{x})}{Z}\right) \geq 0$$

The loss function is lower bounded by the physical free energy (Gibbs-Bogoliubov-Feynman inequality)

# “Reparametrization trick”

How to compute the gradient w.r.t random sampling ?

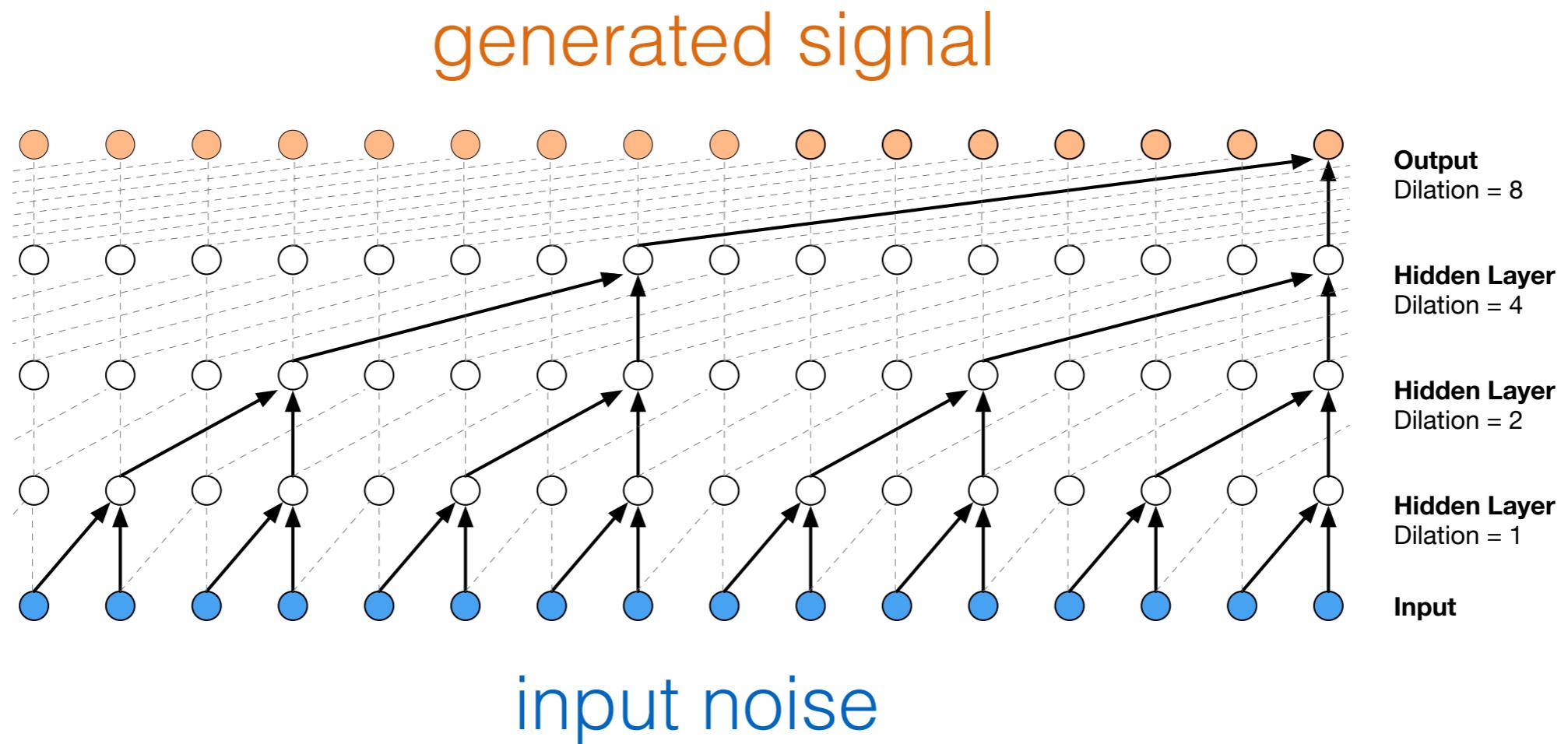
$$\mathcal{L}_\theta = \mathbb{E}_{z \sim p(z)} [\ln q(g_\theta(z)) - \ln \pi(g_\theta(z))]$$

Sample from the prior distribution

Network parameters

End-to-end training using **back-propagation**

# Interlude: WaveNet Story



Given a parallel WaveNet student  $p_S(x)$  and WaveNet teacher  $p_T(x)$  which has been trained on a dataset of audio, we define the **Probability Density Distillation** loss as follows:

$$D_{\text{KL}}(P_S || P_T) = H(P_S, P_T) - H(P_S) \quad (6)$$

# Demo: Ising model

$$\pi(s) = \exp\left(\frac{1}{2}s^T K s\right)$$

# Demo: Ising model

$$\pi(\mathbf{s}) = \exp\left(\frac{1}{2}\mathbf{s}^T K \mathbf{s}\right)$$

decouple

$$\propto \int d\mathbf{x} \exp\left(-\frac{1}{2}\mathbf{x}^T (K + \alpha I)^{-1} \mathbf{x} + \mathbf{s}^T \mathbf{x}\right)$$

# Demo: Ising model

$$\pi(\mathbf{s}) = \exp\left(\frac{1}{2}\mathbf{s}^T K \mathbf{s}\right)$$

decouple

$$\propto \int d\mathbf{x} \exp\left(-\frac{1}{2}\mathbf{x}^T (K + \alpha I)^{-1} \mathbf{x} + \mathbf{s}^T \mathbf{x}\right)$$

trace out  $\mathbf{s}$

$$\boxed{\pi(\mathbf{x}) = \exp\left(-\frac{1}{2}\mathbf{x}^T (K + \alpha I)^{-1} \mathbf{x}\right) \prod_i \cosh(x_i)}$$

# Demo: Ising model

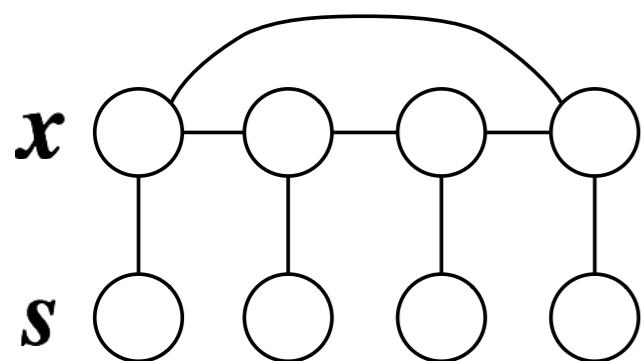
$$\pi(\mathbf{s}) = \exp\left(\frac{1}{2}\mathbf{s}^T K \mathbf{s}\right)$$

decouple

$$\propto \int d\mathbf{x} \exp\left(-\frac{1}{2}\mathbf{x}^T (K + \alpha I)^{-1} \mathbf{x} + \mathbf{s}^T \mathbf{x}\right)$$

trace out  $\mathbf{s}$

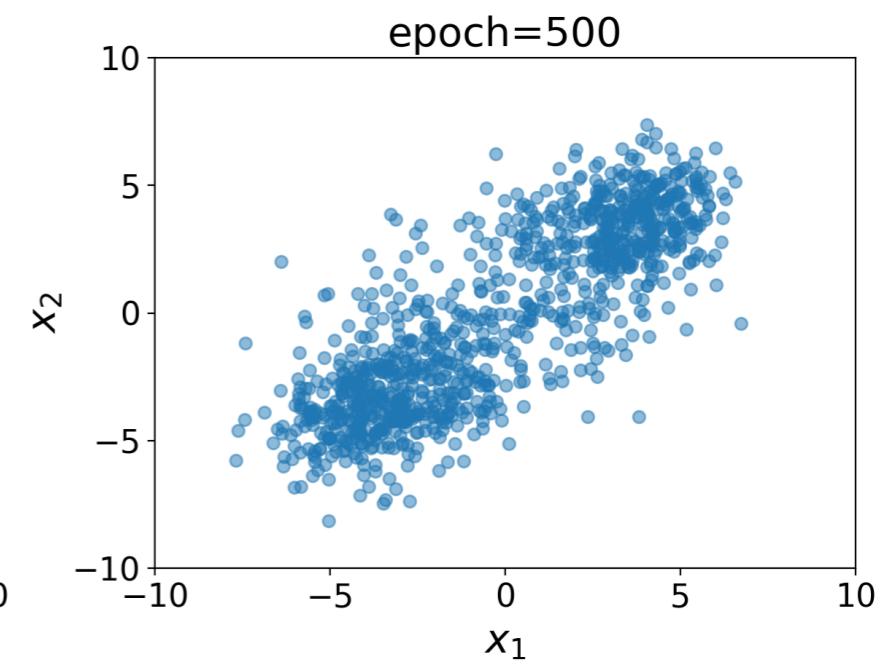
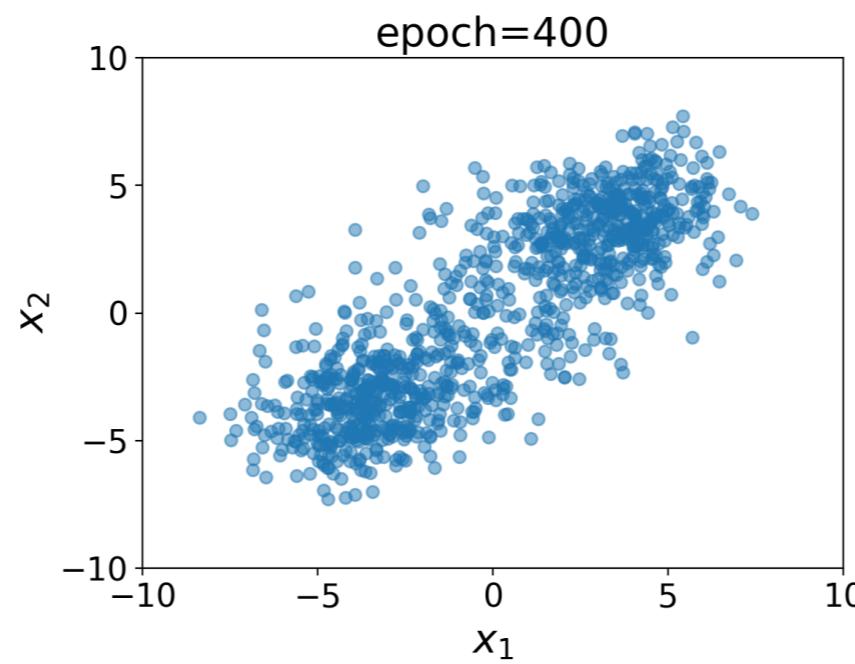
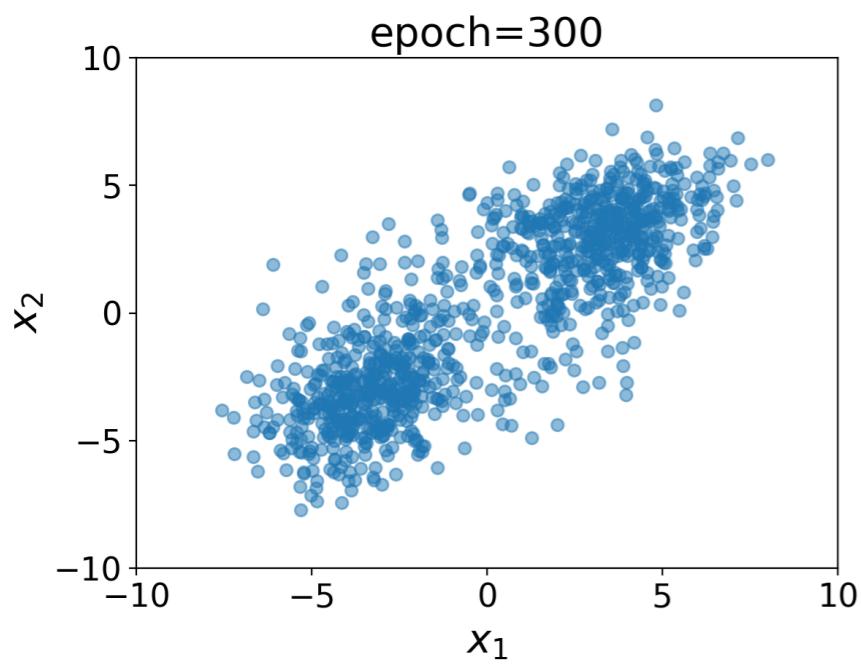
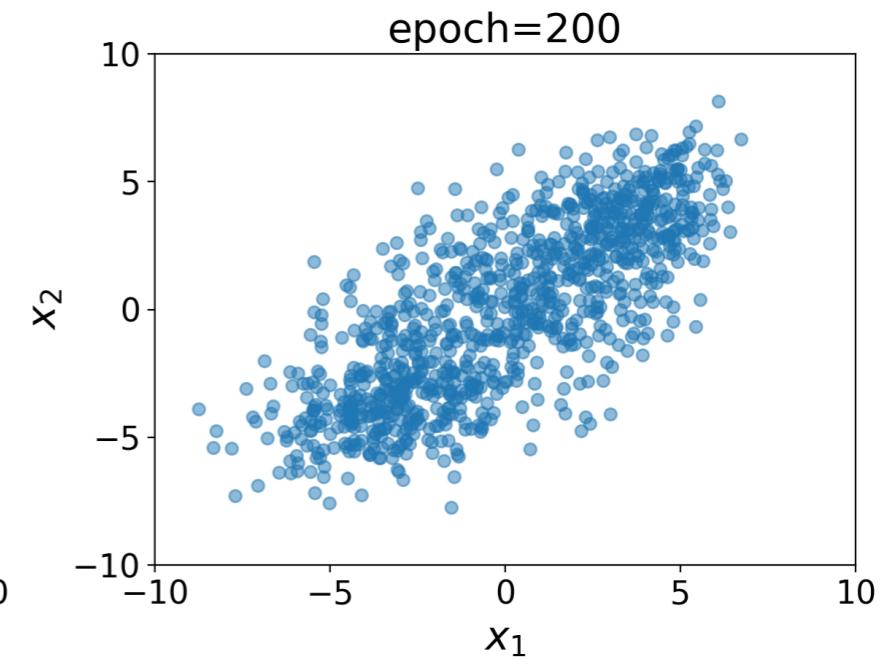
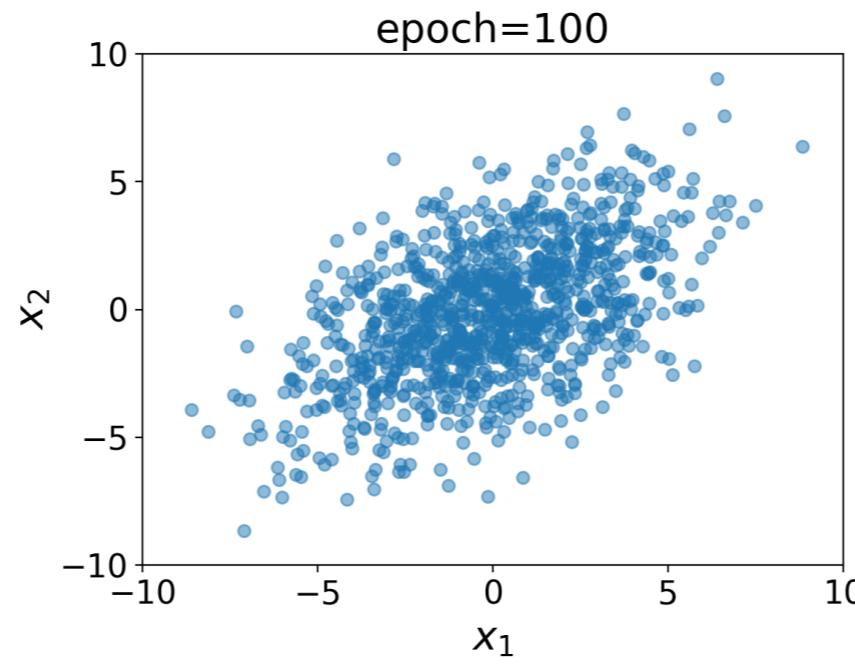
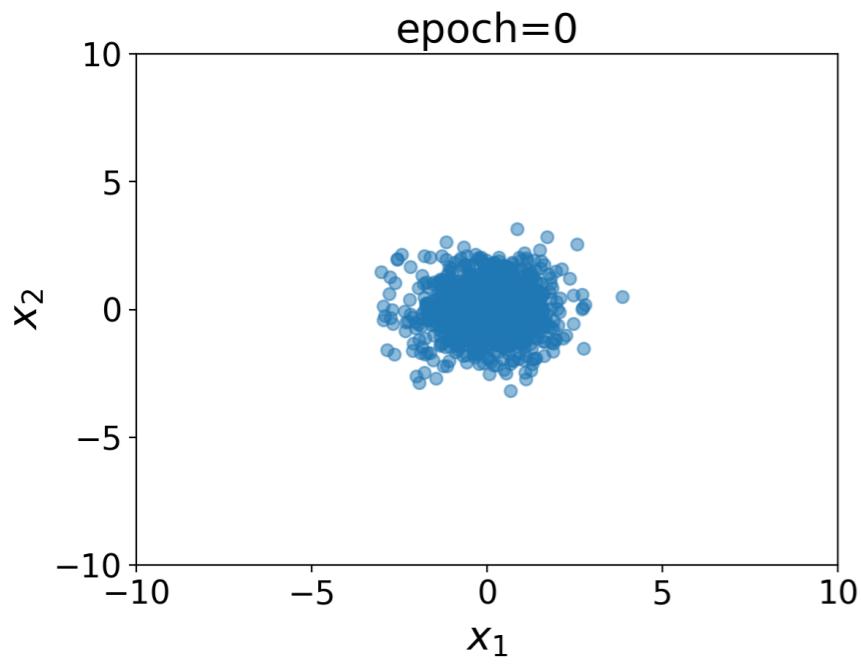
$$\pi(\mathbf{x}) = \exp\left(-\frac{1}{2}\mathbf{x}^T (K + \alpha I)^{-1} \mathbf{x}\right) \prod_i \cosh(x_i)$$



$$\pi(\mathbf{s}|\mathbf{x}) = \prod_i \left(1 + e^{-2s_i x_i}\right)^{-1}$$

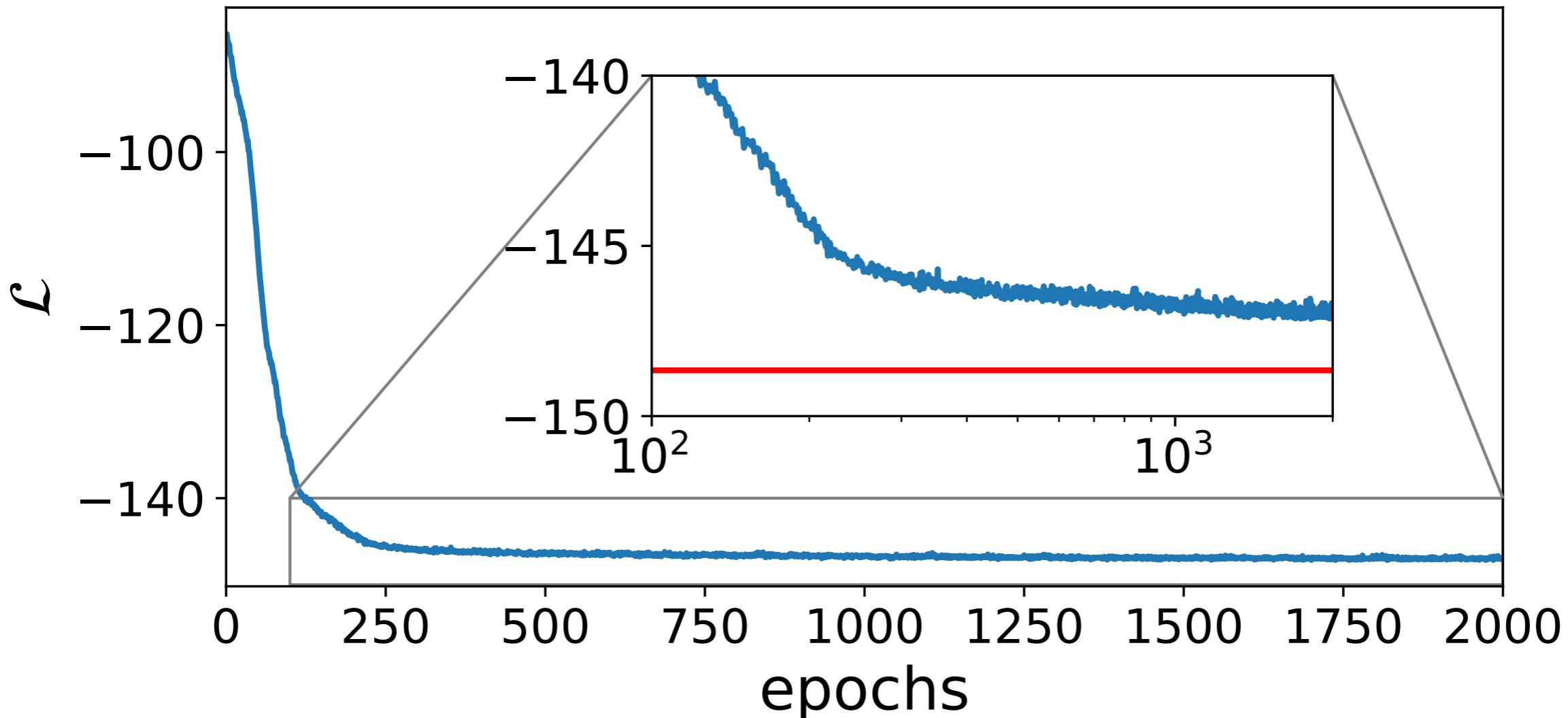
continuous dual  
of the Ising model

# Generated Samples



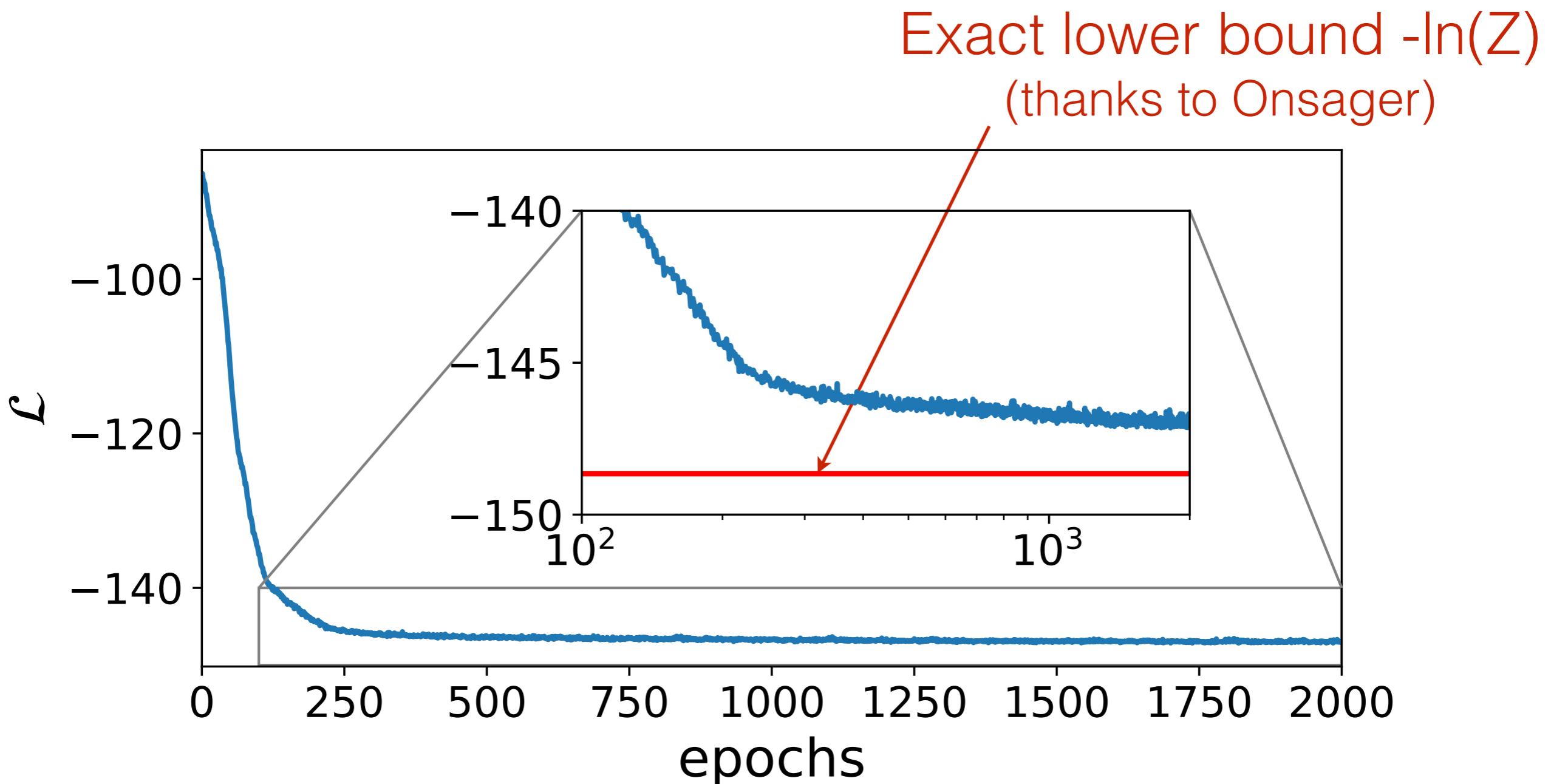
$x$ 's are continuous fields dual to Ising spins

# Variational Loss



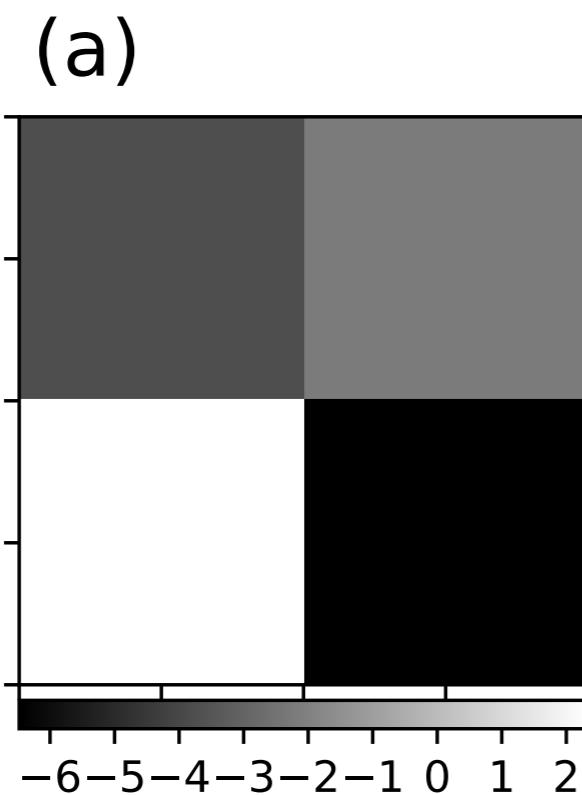
Loss can be further improved by  
using more expressive networks

# Variational Loss

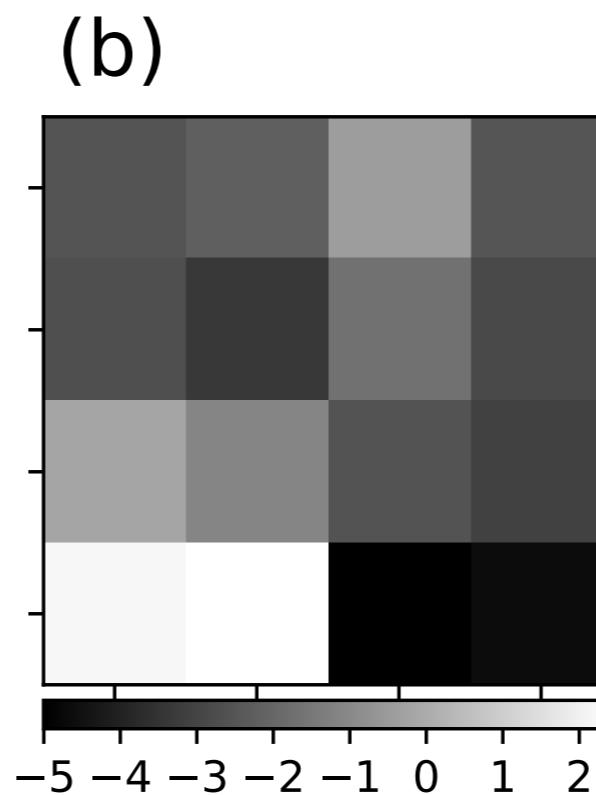


Loss can be further improved by  
using more expressive networks

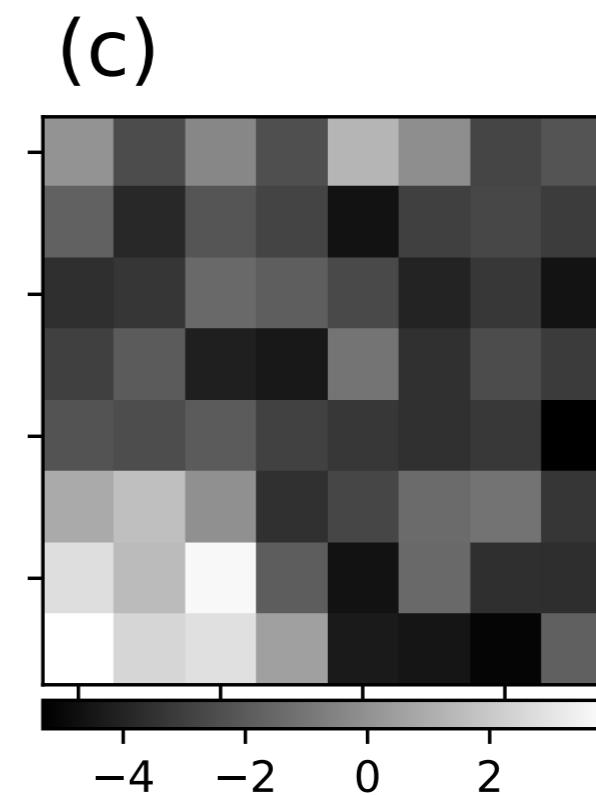
# Renormalized Collective Variables



2x2



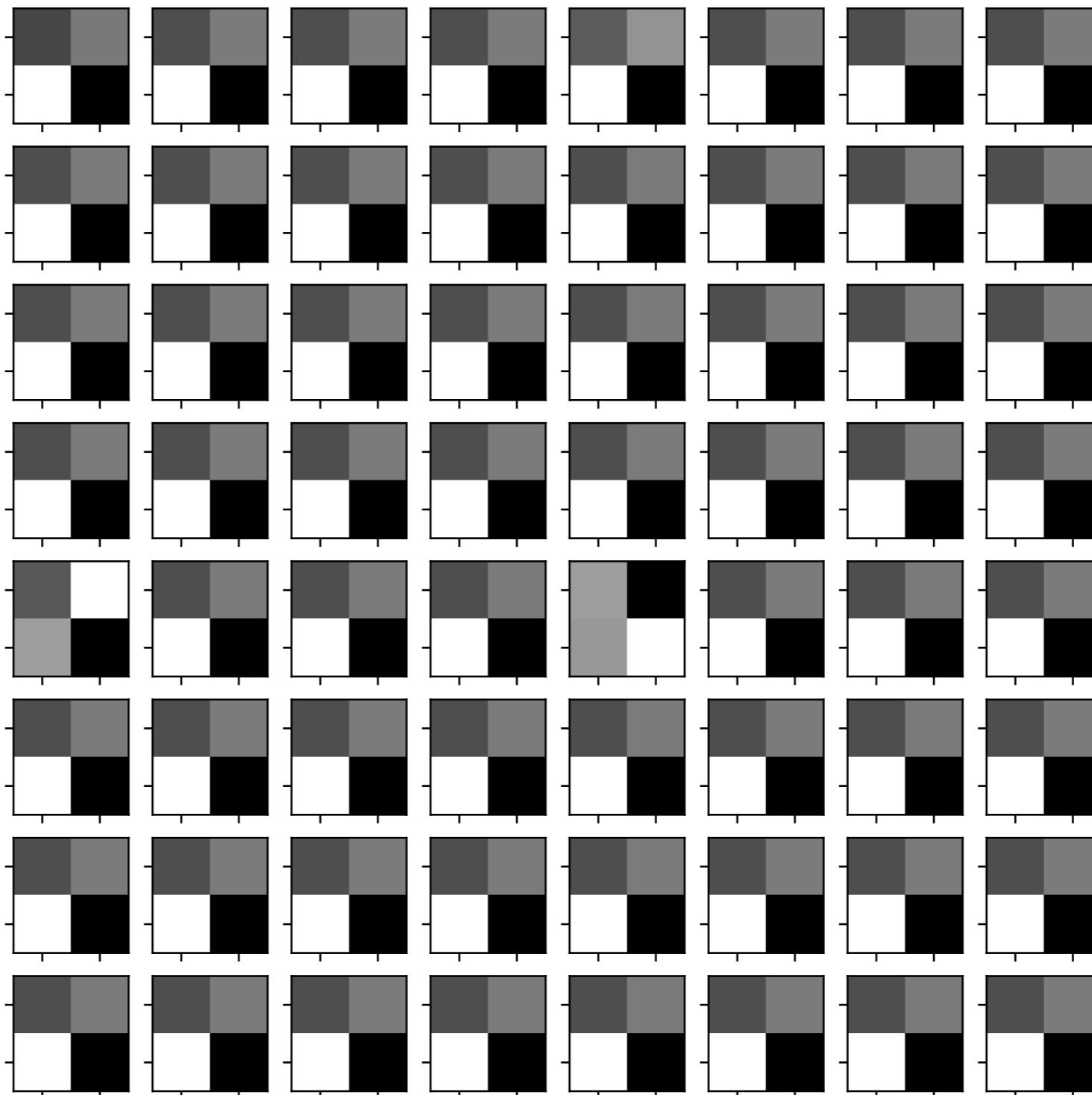
4x4



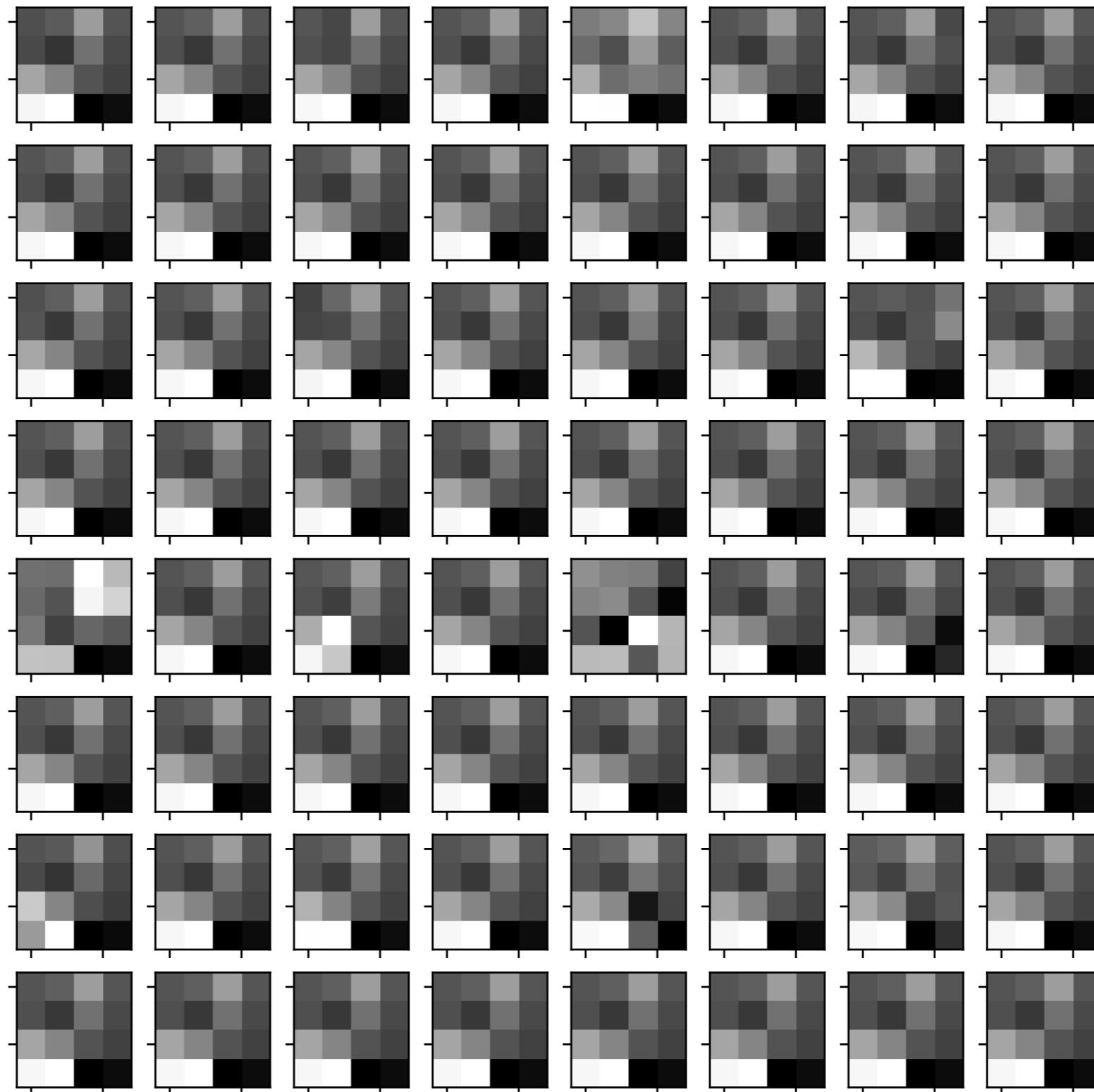
Physical  
Variables

Also know their effective couplings  
=> renormalized energy function

# Wander in the Latent Space



# Wander in the Latent Space

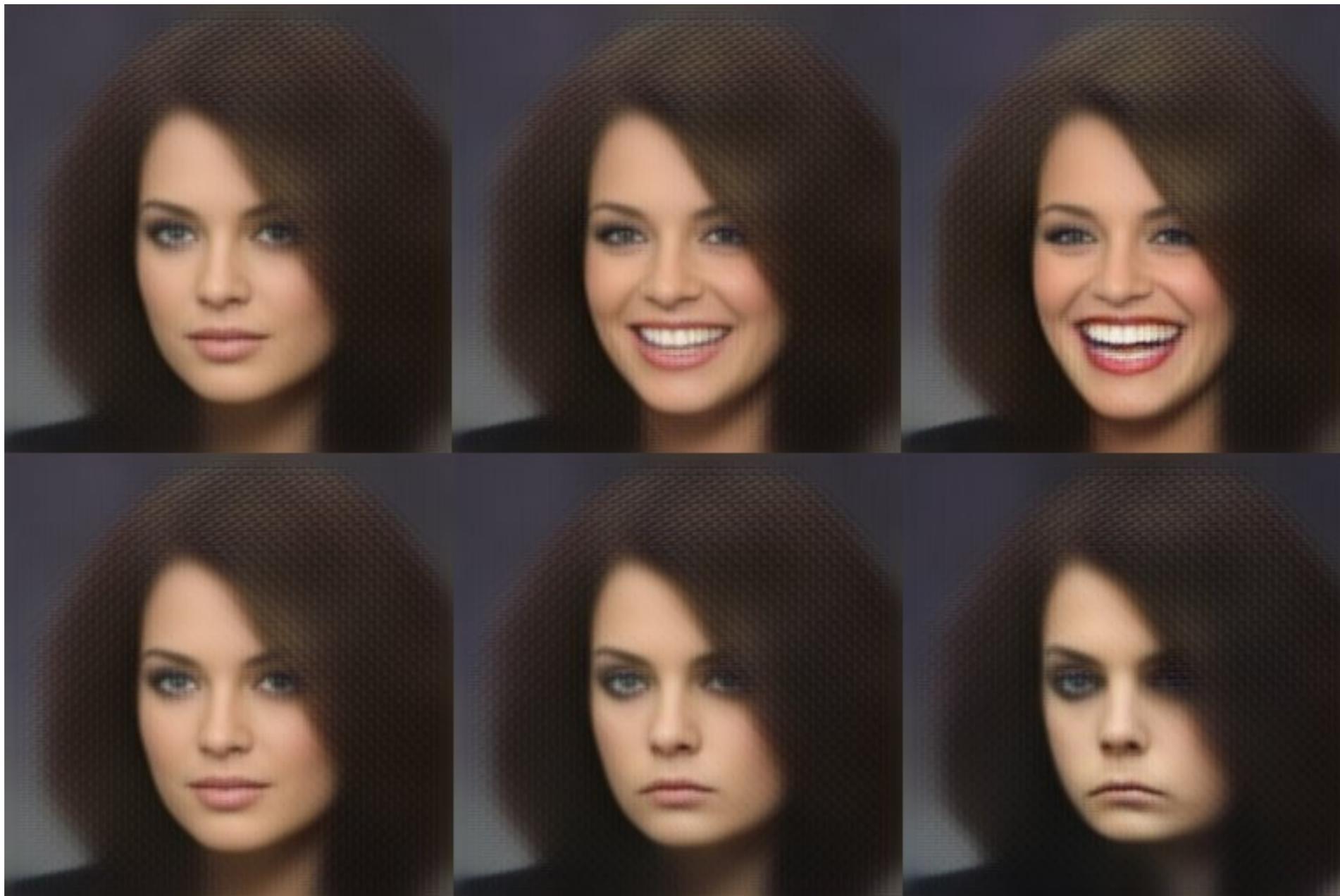


# Wander in the Latent Space



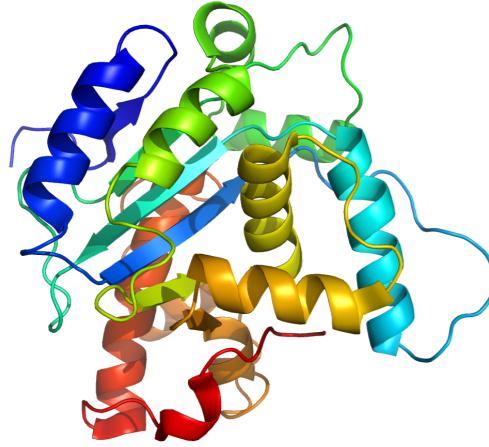
# Wander in the Latent Space

**Arithmetics of the “smile vector”**

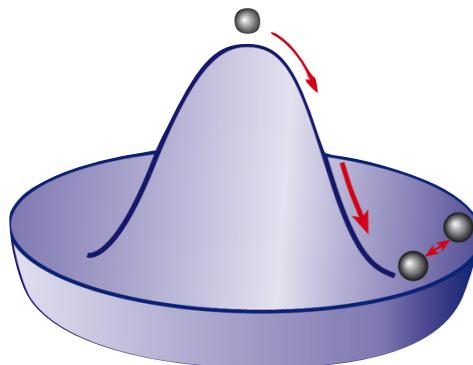


White, 1609.04468 implemented using the variational autoencoder by Kingma and Welling, 1312.6114

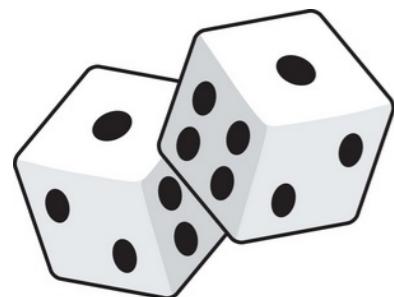
# How is it useful ?



Automatically identify collective variables  
(metadynamics molecular simulation)



Automatically derive effective field theory  
(free energy surface)



Monte Carlo update proposals

# Sampling in the latent space

# Change-of-variables in a learnable way

$$Z = \int dx \pi(x) = \int dz \pi(g(z)) \left| \det \left( \frac{\partial x}{\partial z} \right) \right|$$

↑  
Physical Prob. Dist.

↑  
Latent variable Prob. Dist.

Latent space is less correlated,  
therefore, easier to sample

# Metropolized Independent Sampler

Acceptance rate with detailed balance condition

$$A(\mathbf{x} \rightarrow \mathbf{x}') = \min \left[ 1, \frac{q(\mathbf{x})}{q(\mathbf{x}')} \cdot \frac{\pi(\mathbf{x}')}{\pi(\mathbf{x})} \right]$$

- Unbiased physics even for imperfect proposals
- Proposals are independent

Propose Ratio      Physical Probability

**Surrogate energy function:**

Li Huang and LW, 1610.02746

Liu, Qi, Meng, Fu, 1610.03137

**Trainable transition kernel:**

Song, Zhao, Ermon, 1706.07561

Levy, Hoffman, Sohl-Dickstein, 1711.09268

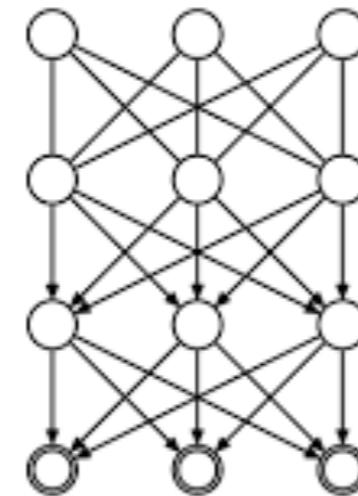
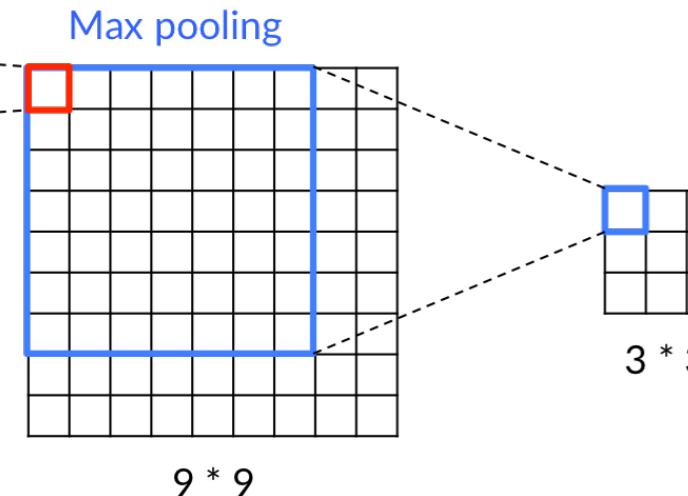
# Remarks on TNS Connection

- What we had is a **classical downgrade of MERA** (Bény 2013)
  - Probability Density ~ Quantum Wavefunction
  - Classical Mutual Information ~ Entanglement Entropy
  - “Decorrelator” ~ Disentangler
  - Decimator ~ Isometry
  - Bijector ~ Unitary
- Deep Learning machinery provides **structural flexibility, modular abstraction, end-to-end training**
- We give back to DL **understandings of what are they doing** (and hopefully, how to do better)

# Remarks on DL

## Old Wisdoms

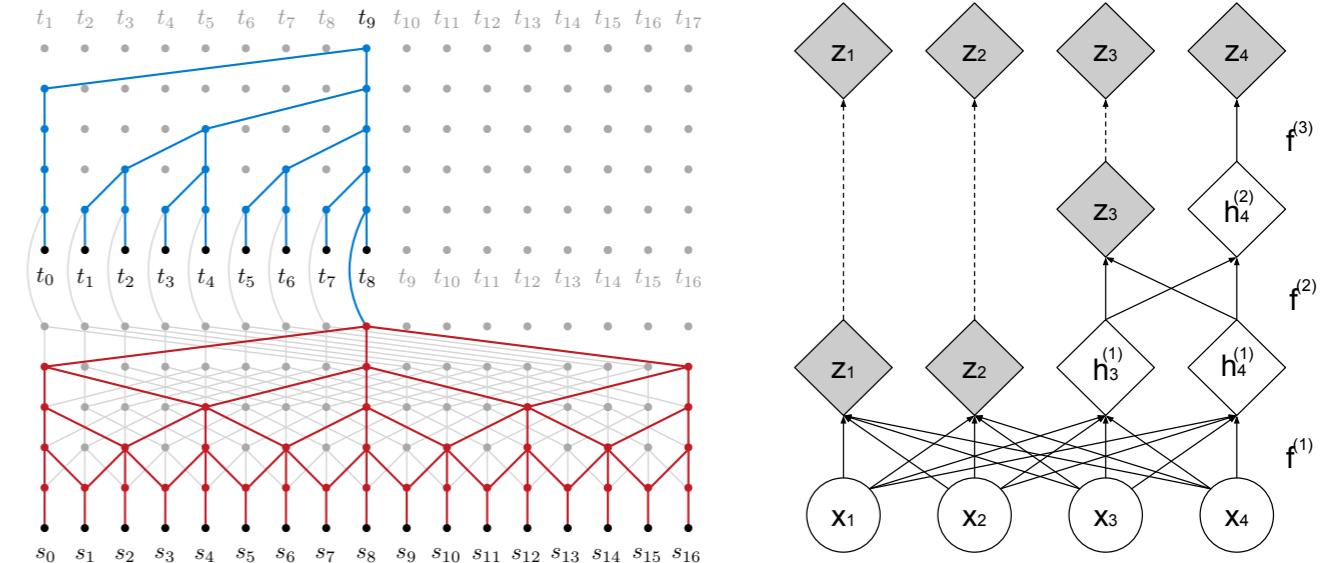
Pooling layer in ConvNets  
~ Decimation



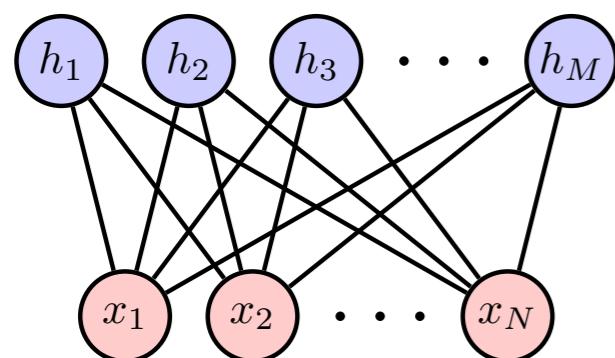
## New Insights

Dialed convolution + Factor out layers = Decimation

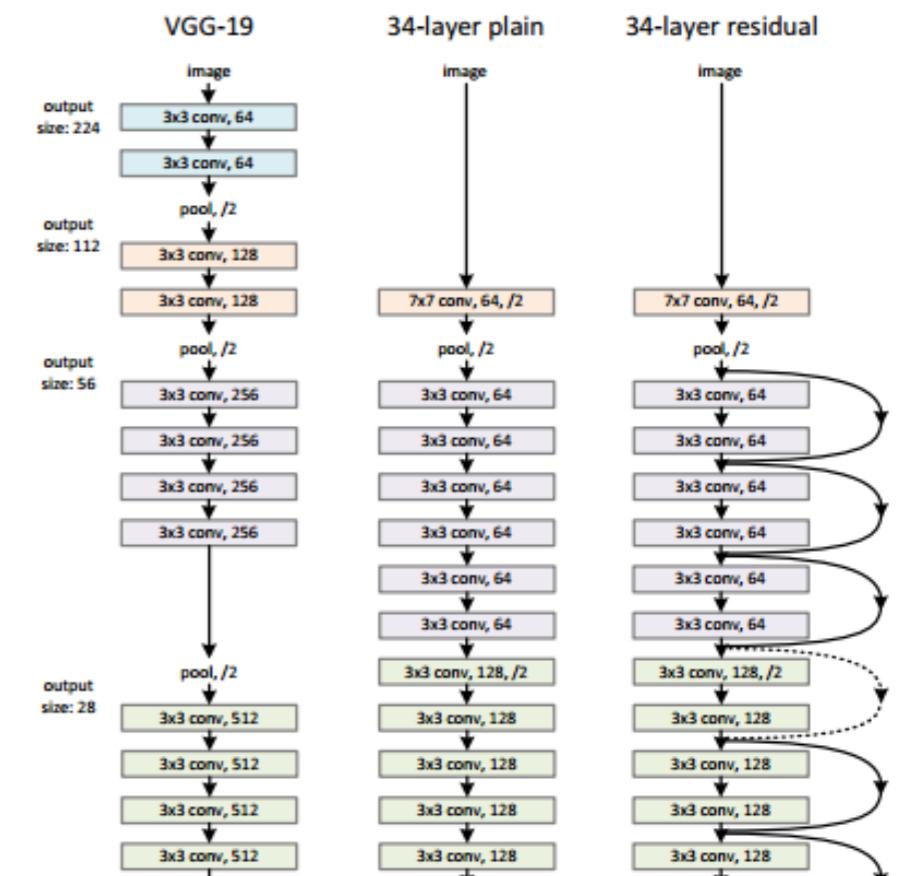
Kept latent variables =  
Renormalized Variables



# Spherical chicken in vacuum



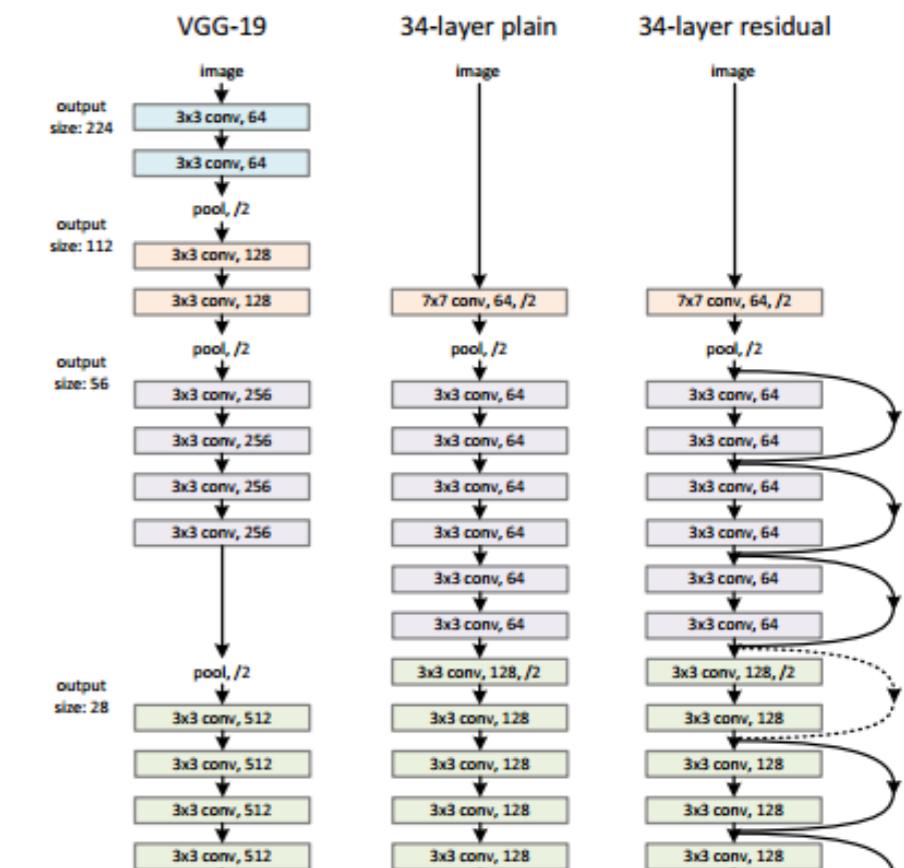
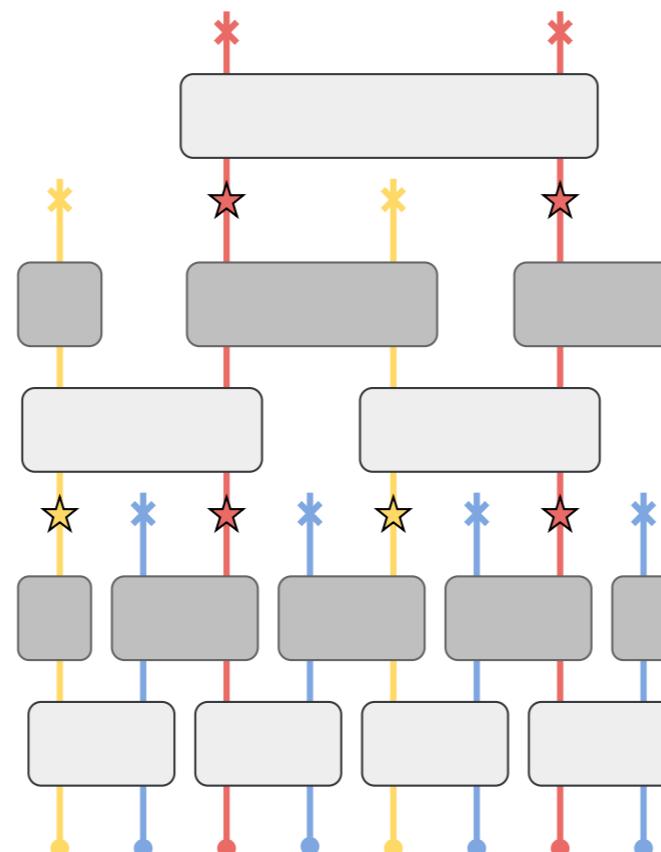
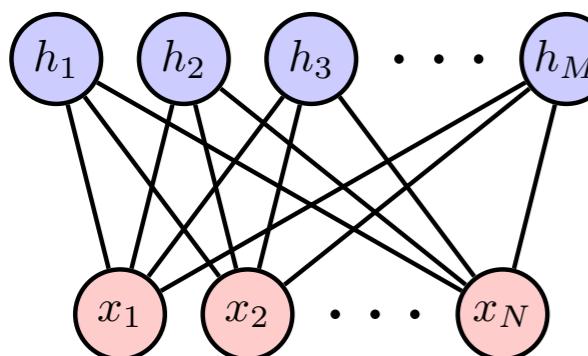
# Animals in the wild



# Here we are

Spherical chicken  
in vacuum

Animals  
in the wild



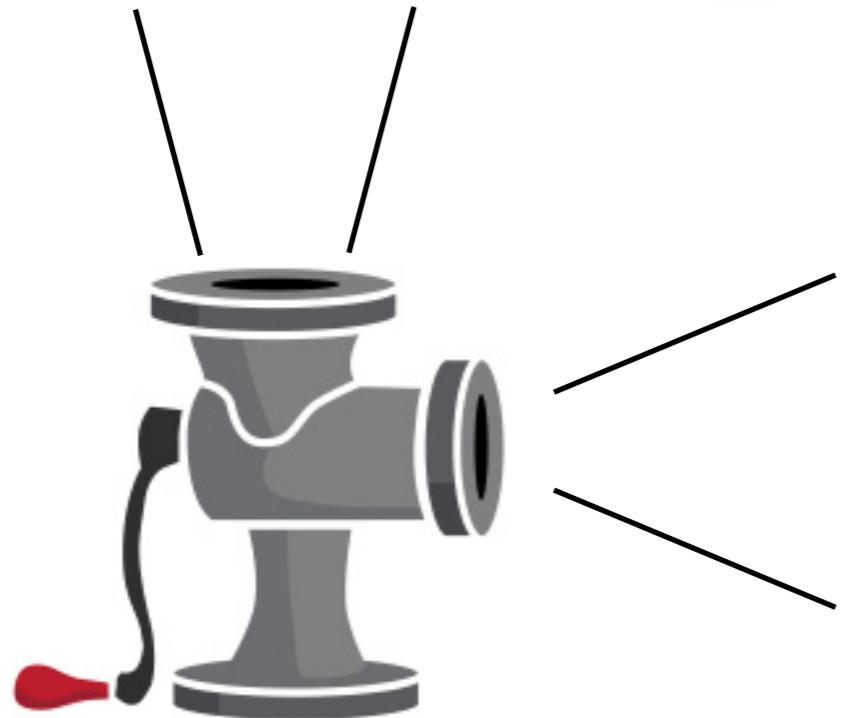
Simplified, but not oversimplified model with balanced interpretability and expressibility

# Remarks on RG

- Conventionally, RG is a **semi-group**, not a **group**
- NeuralRG builds on bijectors, hence a group  
(coarse-graining due to the multiscale structure)
- **Probabilistic** (Jona-Lasinio 75') and **Information Theory** (Apenko 09') views on RG (same is true for neural & tensor networks)
- Diffeomorphism does not change **topology** of the manifolds, therefore, may be limited.

# The Universe as a Generative Model

$$\mathcal{L} = \int d\bar{x} \sqrt{-g} \left[ \frac{m_p^2}{2} R - \frac{1}{4} F_{\mu\nu}^a F_a^{\mu\nu} + i \bar{\psi}^i \gamma^\mu D_\mu \psi^i + \left( \bar{\psi}_L^i \gamma_j \not{D} \psi_R^j + \text{h.c.} \right) - |\partial_\mu \not{\Phi}|^2 - V(\not{\Phi}) \right]$$



**Thank you!**